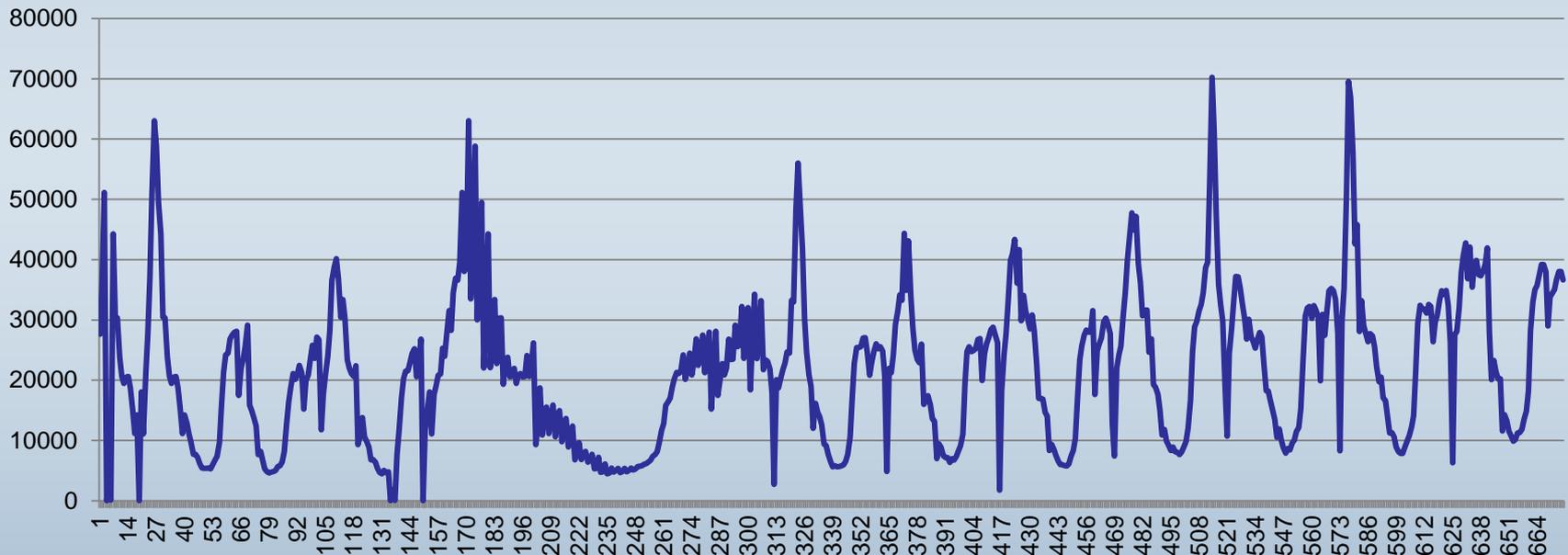


Методы и модели анализа прогнозирования.



Кольцов С.Н

Основные понятия и определения

Два метода прогнозирования: 1. Количественное прогнозирование.
2. Качественное прогнозирование (опрос экспертов).

Сложности прогнозирования:

1. **Очень сложно сделать правдивый прогноз, особенно если он касается будущего.**
2. **Предсказать не сложно, сложно предсказать правильно 😊**
3. **Числа, если их правильно преподнести, могут сказать о чем угодно.**

Количественное прогнозирование: 1. Каузальные модели (причинно – следственные).
2. Временные модели (например, модель Брауна).

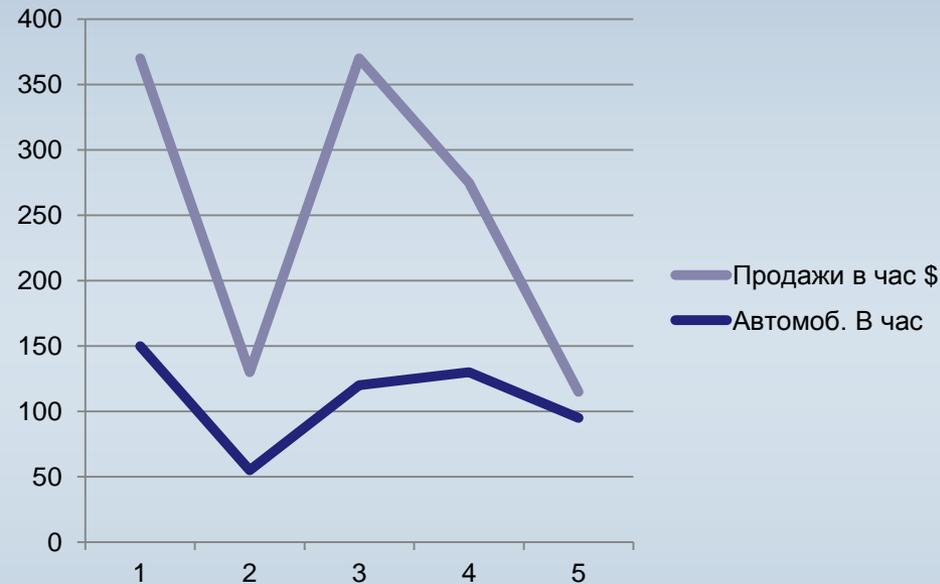
Математическая постановка каузальной задачи $y = f(x_1, x_2, x_3...)$

Например: y – спрос на детскую еду, x – число детей, 1,2,3. – нумерация месяцев (лет)

или: y – спрос на медицинское оборудование, x – число больниц или поликлиник

Причинно – следственная модель

Причинно - следственные модели используются в том случае, когда независимые переменные (x) известны заранее или их спрогнозировать проще чем зависимую переменную y . Например, можно спрогнозировать объем продаж на следующий месяц как функцию от производства за прошлый месяц. Таким образом для выбора П-С модели нужно что бы выполнялись два условия.



1. Должна существовать связь (за прошлый месяц) между независимой переменной и зависимой переменной.
 2. Значения независимой переменной (x) должны быть известны на период, на который мы собираемся делать прогноз.
- Формирование математической связи между зависимой и независимой переменными основано на применение метода наименьших квадратов.

Оценка прогнозирования

Для оценки точности модели можно относительную ошибку аппроксимации и средний квадрат ошибки.

Относительная ошибка

$$error = y - \hat{y},$$

где

y – реальное значение

\hat{y}

y – аппроксим значение

Средний квадрат ошибки

$$error = \frac{(y - \hat{y})^2}{y^2},$$

где

y – реальное значение

\hat{y}

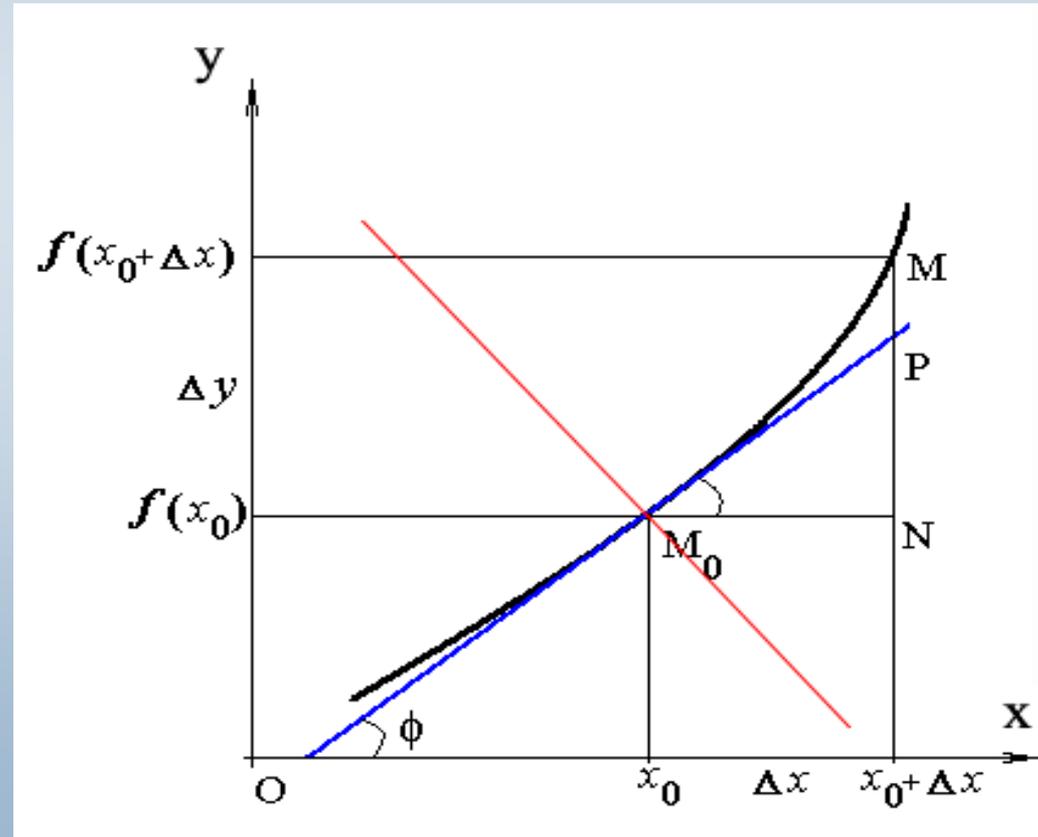
y – аппроксим значение

Считается, что точность модели хорошая, если среднее значение относительной погрешности не превышает 5% , удовлетворительная, если среднее значение относительной погрешности не превышает 15%, и неудовлетворительная, если среднее значение относительной погрешности больше 15%.

Определение и геометрическая интерпретация производной

Пусть x - независимая переменная, $y = f(x)$

Значение производной $f'(x_0)$ равняется угловому коэффициенту касательной к графику функции $y = f(x)$ в точке $M_0(x_0, f(x_0))$; $f'(x_0) = \operatorname{tg} \Phi$, где Φ - угол наклона касательной к оси Ox



$$f'(x_0) = \lim_{\Delta x \rightarrow 0} \frac{f(x_0 + \Delta x) - f(x_0)}{\Delta x}$$

Таблица производных

$C' = 0$	$(\arcsin x)' = \frac{1}{\sqrt{1-x^2}}$
$x' = 1$	$(\arccos x)' = -\frac{1}{\sqrt{1-x^2}}$
$(x^n)' = n \cdot x^{n-1}$	$(\operatorname{arctg} x)' = \frac{1}{1+x^2}$
$(\sqrt{x})' = \frac{1}{2\sqrt{x}}$	$(\operatorname{arcctg} x)' = -\frac{1}{1+x^2}$
$(e^x)' = e^x$	$(\operatorname{sh} x)' = \operatorname{ch} x$
$(a^x)' = a^x \ln a$	$(\operatorname{ch} x)' = \operatorname{sh} x$
$(\ln x)' = \frac{1}{x}$	$(\operatorname{th} x)' = \frac{1}{\operatorname{ch}^2 x}$
$(\log_a x)' = \frac{1}{x \ln a}$	$(\operatorname{cth} x)' = -\frac{1}{\operatorname{sh}^2 x}$
$(\sin x)' = \cos x$	$(\operatorname{arcsh} x)' = \frac{1}{\sqrt{x^2+1}}$
$(\cos x)' = -\sin x$	$(\operatorname{arcch} x)' = \frac{1}{\sqrt{x^2-1}}$
$(\operatorname{tg} x)' = \frac{1}{\cos^2 x}$	$(\operatorname{arcth} x)' = \frac{1}{1-x^2}$
$(\operatorname{ctg} x)' = -\frac{1}{\sin^2 x}$	$(\operatorname{arccth} x)' = \frac{1}{1-x^2}$

Производные простых функций (x – независимая переменная)	Производные сложных функций ($u = u(x)$ – любая дифференцируемая функция)
1. $(x^2)' = 2x$	1. $(u^2)' = 2u \cdot u'$
2. $(\sqrt{x})' = \frac{1}{2\sqrt{x}}$	2. $(\sqrt{u})' = \frac{1}{2\sqrt{u}} \cdot u'$
3. $\left(\frac{1}{x}\right)' = -\frac{1}{x^2}$	3. $\left(\frac{1}{u}\right)' = -\frac{1}{u^2} \cdot u'$
4. $(e^x)' = e^x$	4. $(e^u)' = e^u \cdot u'$
5. $(\ln x)' = \frac{1}{x}$	5. $(\ln u)' = \frac{1}{u} \cdot u'$
6. $(\sin x)' = \cos x$	6. $(\sin u)' = \cos u \cdot u'$
(.....)	(.....)

Метод наименьших квадратов

Если некоторая величина зависит от другой величины, то эту зависимость можно исследовать, измеряя y при различных значениях x . В результате измерений получается ряд значений:

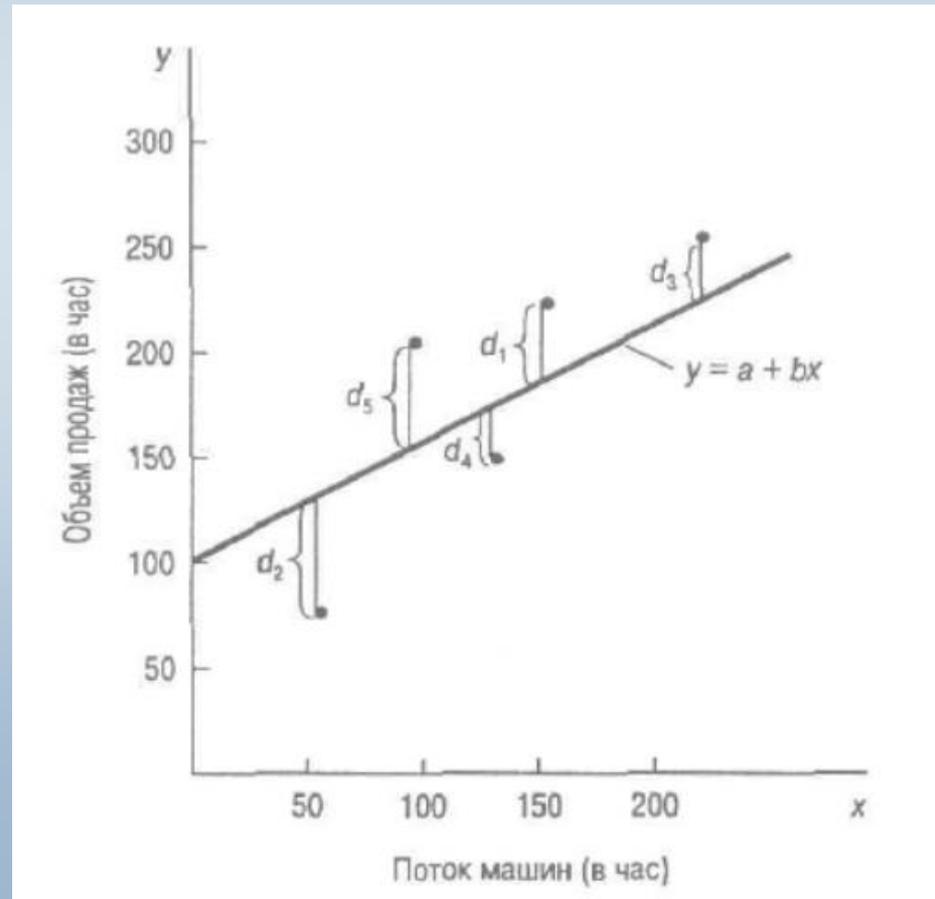
$$x_1, x_2, \dots, x_i, \dots, x_n;$$

$$y_1, y_2, \dots, y_i, \dots, y_n.$$

По данным такого эксперимента можно построить график зависимости

$y = f(x, a, b, c, \dots)$. Полученная кривая дает возможность судить о виде функции $f(x)$.

Однако постоянные коэффициенты, которые входят в эту функцию, остаются неизвестными.



Метод наименьших квадратов для линейной функции

Аппроксимация наших данных, например линейной функцией можно выразить следующим образом:

$$\sum_{i=1}^n d_i^2 = \sum_{i=1}^n (y_i - (a + bx_i))^2$$

Суть метода НКв заключается в том что можно выбрать такие коэффициенты a, b , так что бы минимизировать квадрат разности между аналитической функцией и экспериментальными данными.

Что бы минимизировать квадрат разности нужно найти частные производные по коэффициентам a, b .

$$\sum_{i=1}^n -2(y_i - (a + bx_i)) = 0$$

$$\sum_{i=1}^n -2x_i(y_i - (a + bx_i)) = 0$$

Метод наименьших квадратов для линейной функции

$$\sum_{i=1}^n -2(y_i - (a + bx_i)) = 0$$



$$a = \frac{1}{n} \sum_{i=1}^n y_i - b \frac{1}{n} \sum_{i=1}^n x_i$$

$$\sum_{i=1}^n -2x_i(y_i - (a + bx_i)) = 0$$



$$b = \frac{\sum_{i=1}^n x_i y_i - \frac{1}{n} \sum_{i=1}^n x_i \sum_{i=1}^n y_i}{\sum_{i=1}^n x_i^2 - \frac{1}{n} \left(\sum_{i=1}^n x_i \right)^2}$$

Метод наименьших квадратов в Экселе

Аппроксимация зависимостей и регрессионный анализ

Интервал абсцисс: Data!\$D\$3:\$D\$52
 Интервал ординат: Data!\$E\$3:\$E\$52
 Интервал вывода решения: Data!\$G\$3

Тип модели

- Экспоненциальная функция
- Степенная функция
- Логарифмическая функция
- Эксп.-степенная функция
- Логистическая функция
- Полином
- Гипербола
- Пользовательская функция
- Кусочно-линейная функция

Параметры полиномов

Минимальная степень: 0 | Степень полинома: 2

Параметры пользовательской функции

Интервал параметров:
 Ячейка формулы:
 Рабочая ячейка:

Метод оптимизации

- Бройдена-Флетчера-Голдфарба-Шанно
- Гаусса-Ньютона

Дополнительно

Доверительная вероятность: 0,95

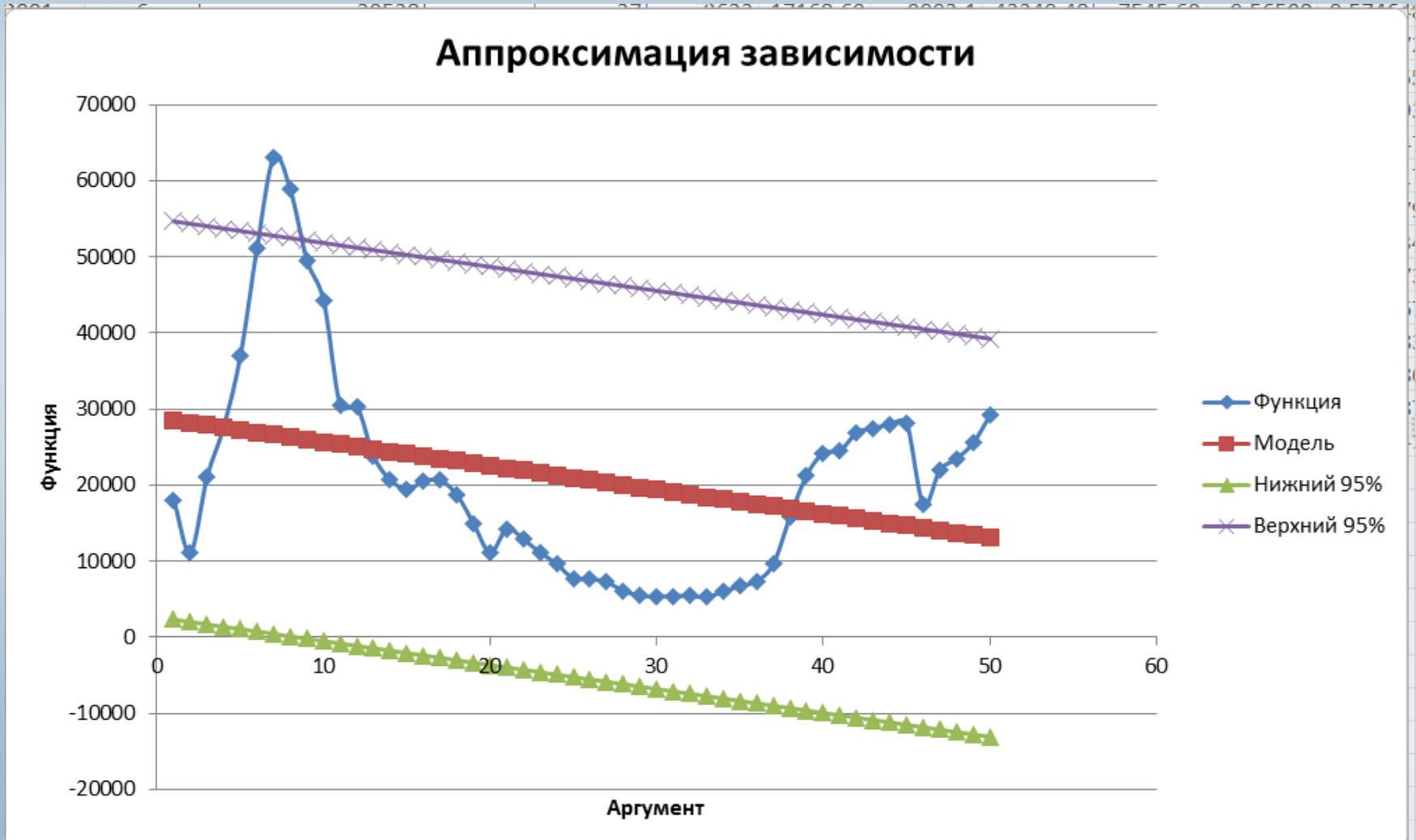
Допустимое число итераций: 100
 Точность: 0,0001

Расчет | Отмена | Помощь

Модуль **APX** - Аппроксимация зависимостей

Справка по APX | Аппроксимация зависимостей

Метод наименьших квадратов в Экселе, линейная функция



Метод наименьших квадратов для полимальной функции

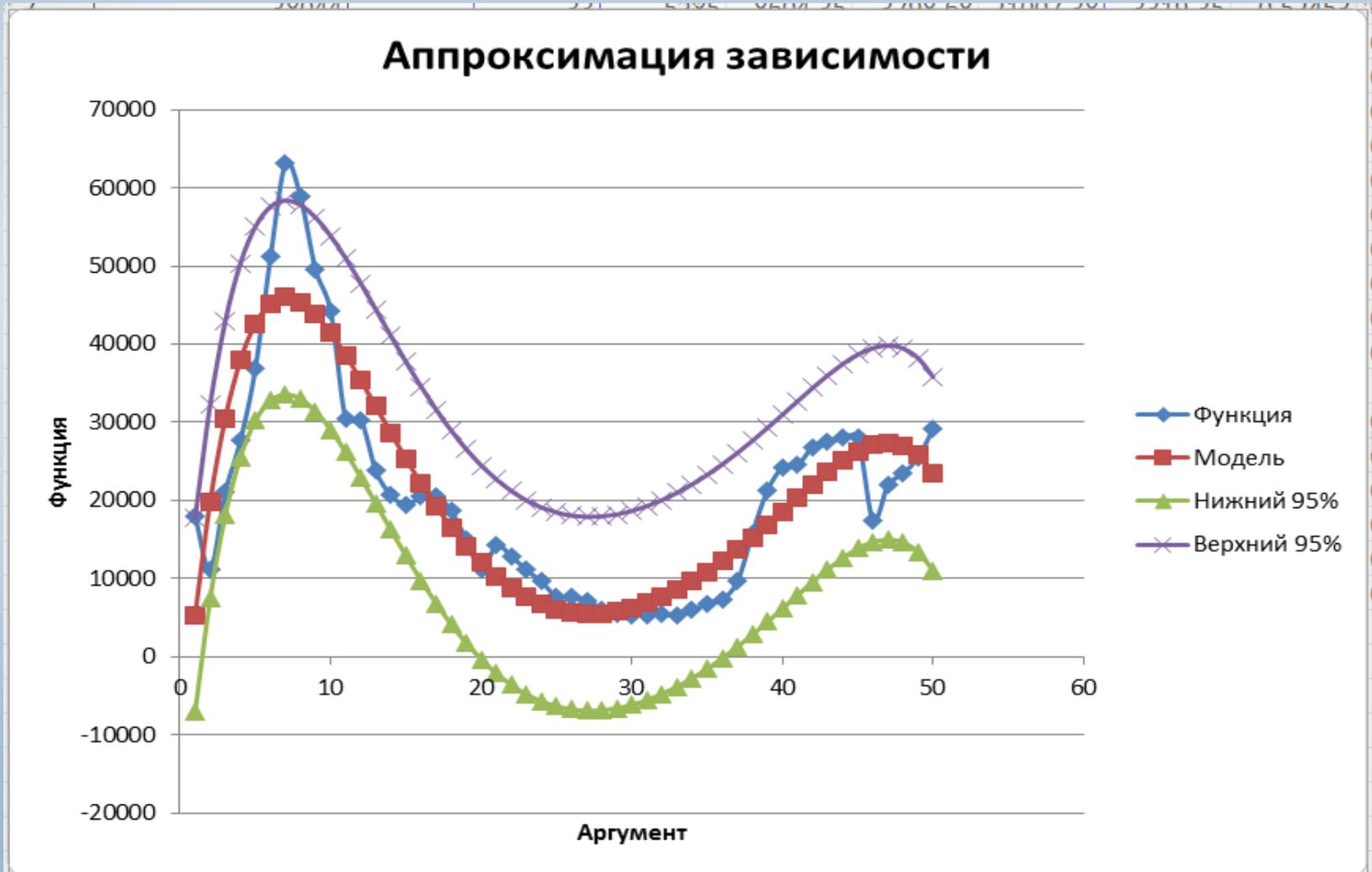
Рассмотрим аппроксимационную функцию в виде квадратичной функции, соответственно ищем квадрат разности между экспериментальными данными и нашей аналитической функцией с не известными коэффициентами.

$$\sum_{i=1}^5 (y_i - (a_0 + a_1 x_i + a_2 x_i^2))^2 .$$

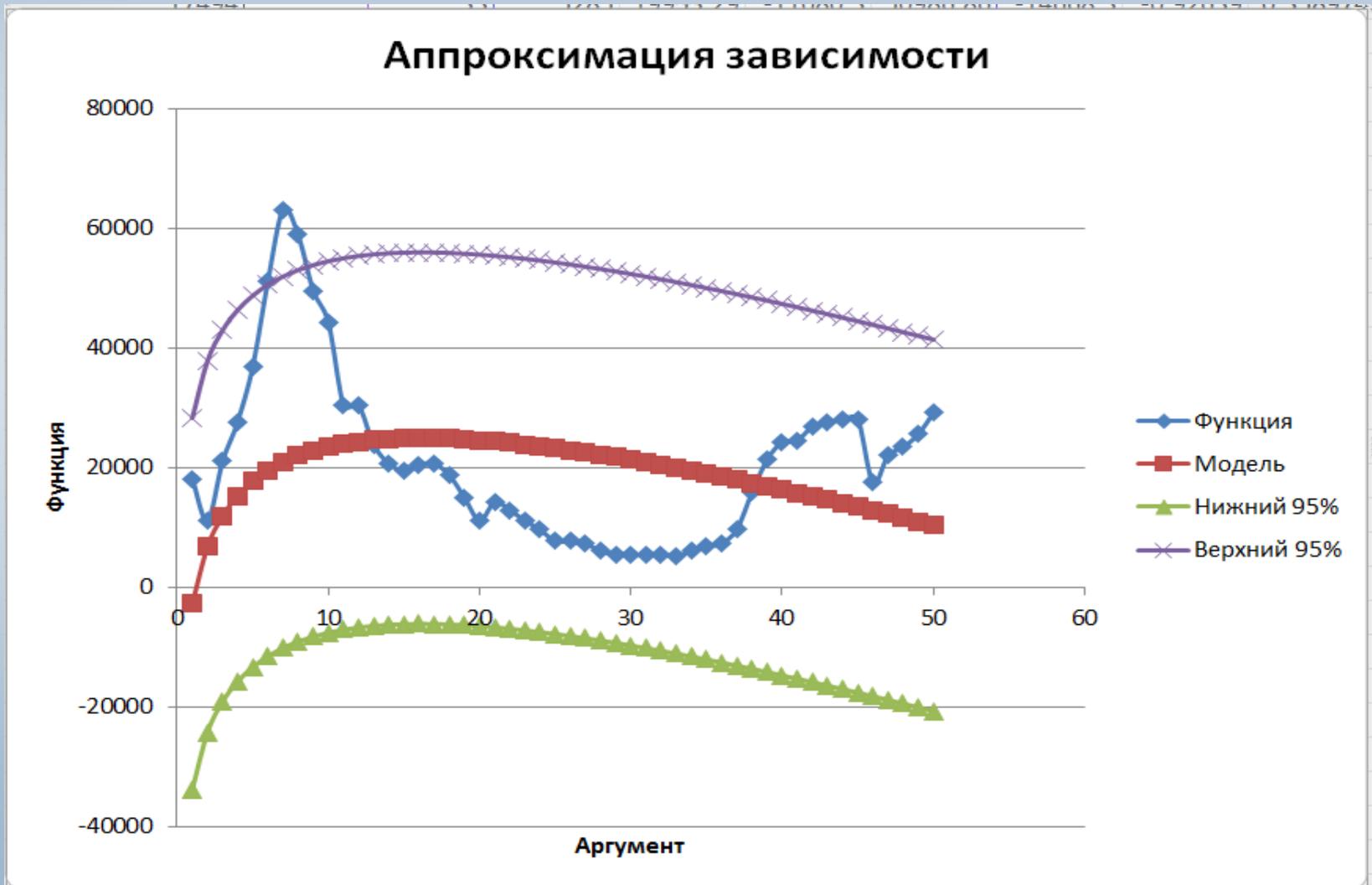
Приравняем нулю три производные по трем параметрам, соответственно получим три набора уравнений:

$$\begin{aligned} 5a_0 + (\sum x_i) a_1 + (\sum x_i^2) a_2 &= \sum y_i , \\ (\sum x_i) a_0 + (\sum x_i^2) a_1 + (\sum x_i^3) a_2 &= \sum x_i y_i , \\ (\sum x_i^2) a_0 + (\sum x_i^3) a_1 + (\sum x_i^4) a_2 &= \sum x_i^2 y_i . \end{aligned}$$

Метод наименьших квадратов в Экселе, полином 6 степени



Метод наименьших квадратов для логарифмической функции.



Пример пользовательской функции

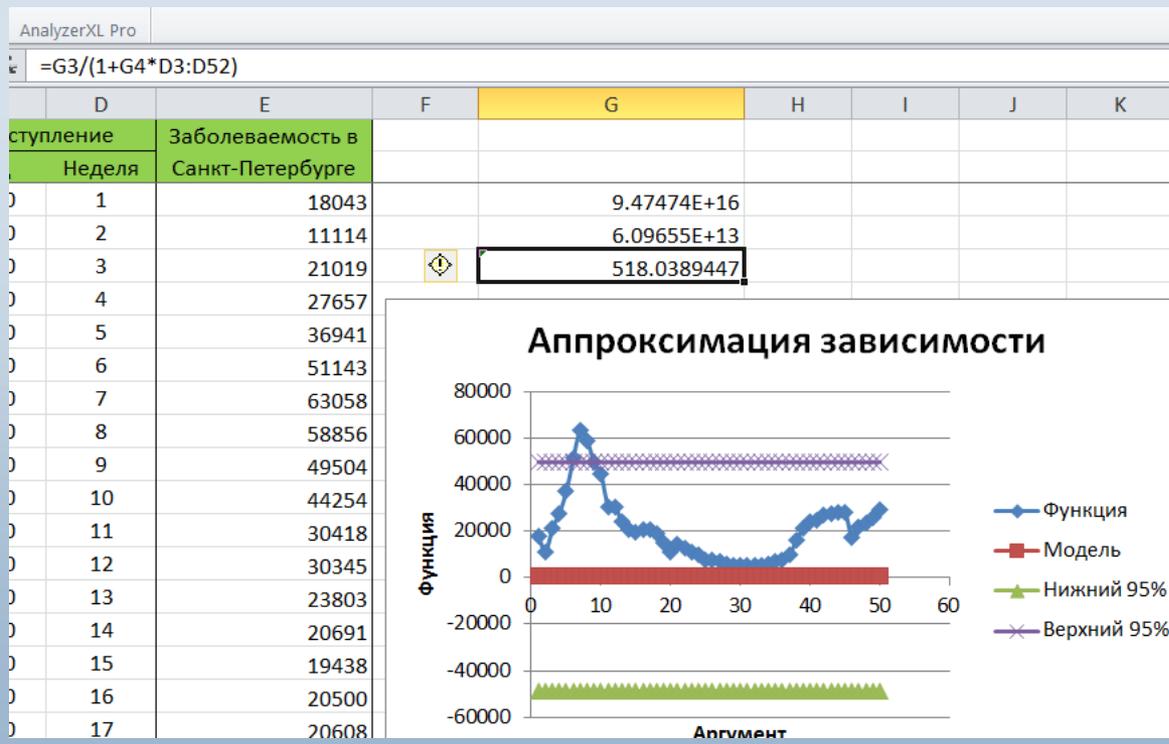
Пусть мы хотим аппроксимировать данные функцией следующего вида:

$$y(t) = \theta_1 / (1 + \theta_2 t)$$

В ходе аппроксимации методом наименьших квадратов ищутся параметры:

$$\theta_1, \theta_2$$

Нам нужно задать начальные данные для этих параметров (в отдельных ячейках), и задать вид функции (также в отдельной ячейке), например



Оценка точности метода наименьших квадратов

Дисперсия функции вычисляется следующим образом:

$$\sigma_Y^2 = \frac{1}{N} \sum_{i=1}^N (y_i - \bar{y})^2$$

где \bar{y} (с чертой) является средним значением, и рассчитывается следующим образом:

$$\bar{y} = \frac{1}{N} \sum_{i=1}^N y_i$$

Метод наименьших квадратов

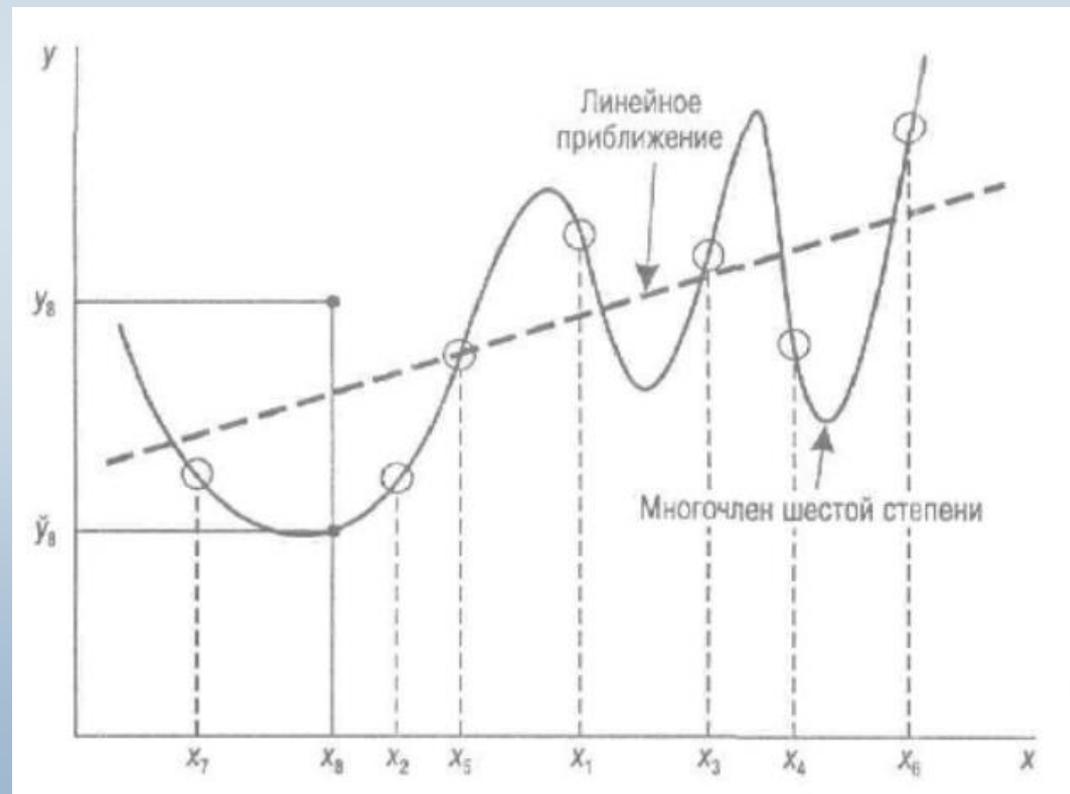
Проблема выбора степени полинома.

Собственно выбирая большую степень полинома мы будем увеличивать точность аппроксимации, то есть будет уменьшаться сумма квадратов.

Однако это приводит к тому что между точками будет возникать большие флуктуации.

Для того что бы выбрать 'наилучшее приближение' можно использовать следующее правило.

Нужно сравнивать среднее значения квадратов ошибок, то, есть для каждого приближения берется величина квадрата ошибки и делится на число точек минус число параметров. Чем меньше эта величина тем лучше.



Модель временных рядов.

В статистике под **временным рядом** понимаются последовательно измеренные через некоторые (зачастую равные) промежутки времени данные. **Анализ временных рядов** объединяет методы изучения временных рядов, как пытающиеся понять природу точек данных, так и пытающиеся построить прогноз. **Прогнозирование временных рядов** заключается в построении модели для предсказания будущих событий основываясь на известных событиях прошлого, предсказания будущих данных до того как они будут измерены. На пример — предсказание цены открытия биржи основываясь на предыдущей её деятельности.

Какие задачи здесь возникают?

1. **Физика солнца:** а) скрытые периодичности; б) прогноз активности.

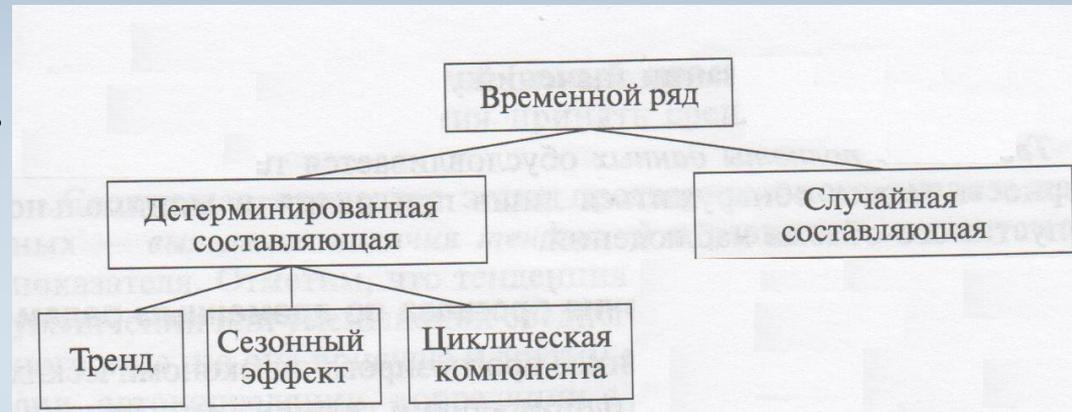
Электрокардиограммы (ЭКГ): а) природа наблюдающихся аритмий;
б) прогноз развития состояния.

Экономические ряды: а) задача сегментации; б) задача прогноза.

Химическая кинетика: а) анализ динамики; б) построение модели.

Модель временных рядов.

Тренд (тенденция) - устойчивая закономерность, наблюдаемую в течении длительного времени. Например, демографическая характеристика или рост потребления.



Сезонная компонента - функция, которая характеризует сезонные колебания. Как правило это колебания носят периодический или почти периодический характер, например загруженность дорог, пик продаж товаров.

Циклическая компонента - функция, описывающая длительные периоды (более года) спада или подъема. примером таких колебаний являются волны Кондратьева, демографические ямы.

Случайная компонента - шум, который отражает случайное действие многочисленных факторов.

Обзор моделей прогнозирования на основе временных рядов.

1. Регрессионные модели прогнозирования:

- Простая линейная регрессия (linear regression)
- Множественная регрессия (multiple regression)
- Нелинейная регрессия (nonlinear regression)

2. Авторегрессионные модели прогнозирования:

ARIMAX (autoregression integrated moving average extended)

GARCH (generalized autoregressive conditional heteroskedasticity),

ARDLM (autoregression distributed lag model)

3. Модели экспоненциального сглаживания (ES):

Взвешенное скользящее среднее

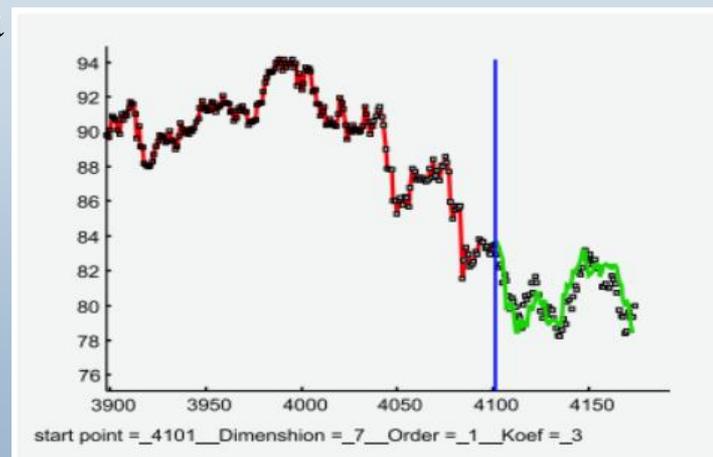
Экспоненциальное сглаживание (exponential smoothing) (Модель Брауна)

Модель Хольта

Модель Хольта-Винтерса

Модель временных рядов.

4. Модель по выборке максимального подобия (MMSP)
5. Модель на нейронных сетях (ANN)
6. Модель на цепях Маркова (Markov chains)
7. Модель на классификационно-регрессионных деревьях (CART)
8. Модель на основе генетического алгоритма (GA)
9. Модель на опорных векторах (SVM)
10. Модель на основе передаточных функций (TF)
11. Модель на нечеткой логике (FL)
12. Модель сингулярного спектрального анализа (SSA)
13. Модель локальной аппроксимации (LA)
14. Модель фрактальных временных рядов
15. Модель на основе Вейвлет - преобразования
16. Модель на основе Фурье-преобразования.



Экспоненциальное сглаживание (Модель Брауна).

Модели сглаживания относятся к адаптивным моделям прогнозирования, способным изменять свою структуру и параметры, приспособившись к изменению условий. Все адаптивные модели делятся на два класса:

1. Модели скользящего среднего (СС-модели) (различные вариации)
2. Авторегрессии (АР-модели). – **это мы не рассматриваем**

Согласно схеме скользящего среднего оценкой текущего уровня (наблюдения) является взвешенное среднее всех предшествующих уровней, причем вес (множитель), который отражает информационную ценность наблюдения, тем больше, чем ближе оно находится к текущему уровню. Такие модели хорошо отражают тенденцию, но не позволяют отражать колебания.

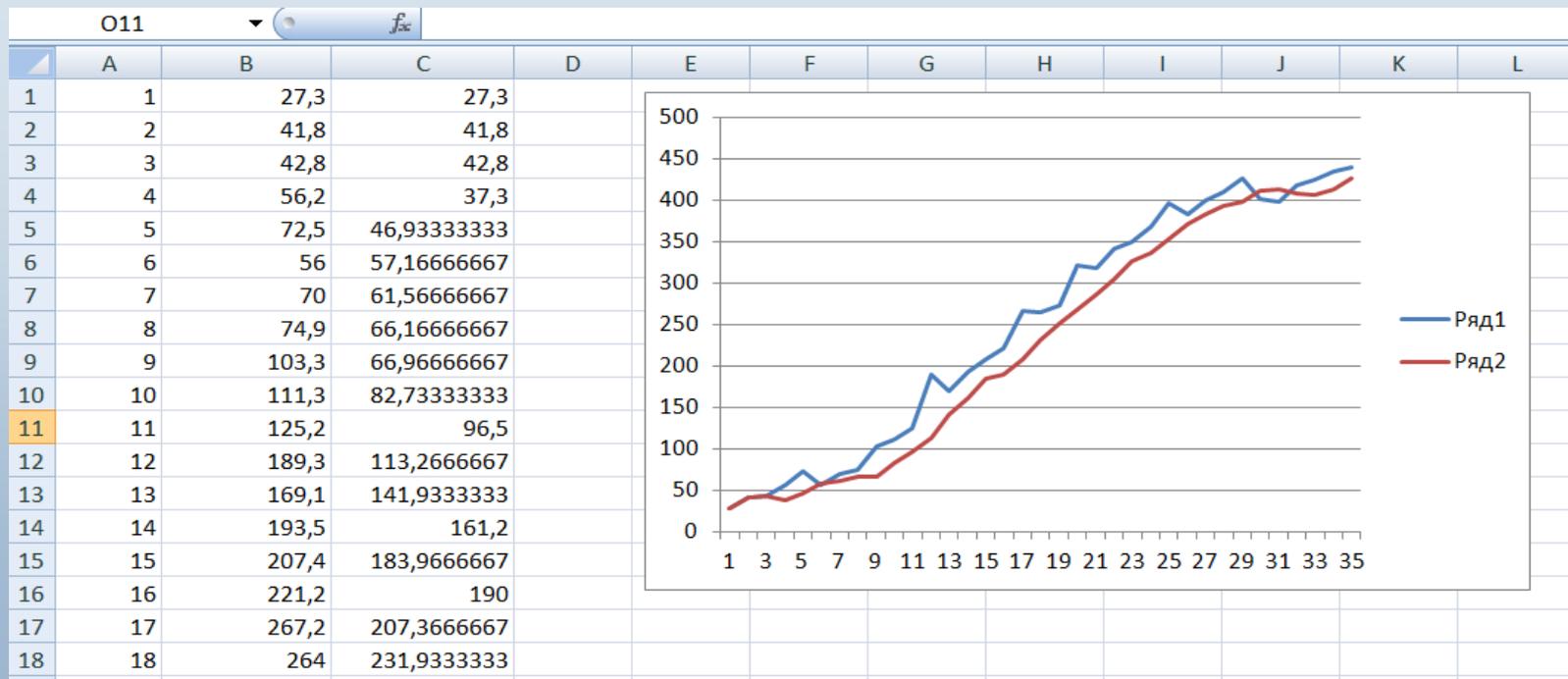
В СС - моделях сглаживание производится с помощью параметра сглаживания, который принимает значения в интервале от 0 до 1.

Скользящее среднее.

Метод скользящего среднего заключается в том что предсказанное значение вычисляется как среднее арифметическое по трем предыдущим данным:

Этот метод годится для медленно меняющихся кривых, с небольшой шумовой компонентой.

$$\hat{y}_{t+1} = \frac{1}{n} (y_t + y_{t-1} + \dots + y_{t-n+1})$$



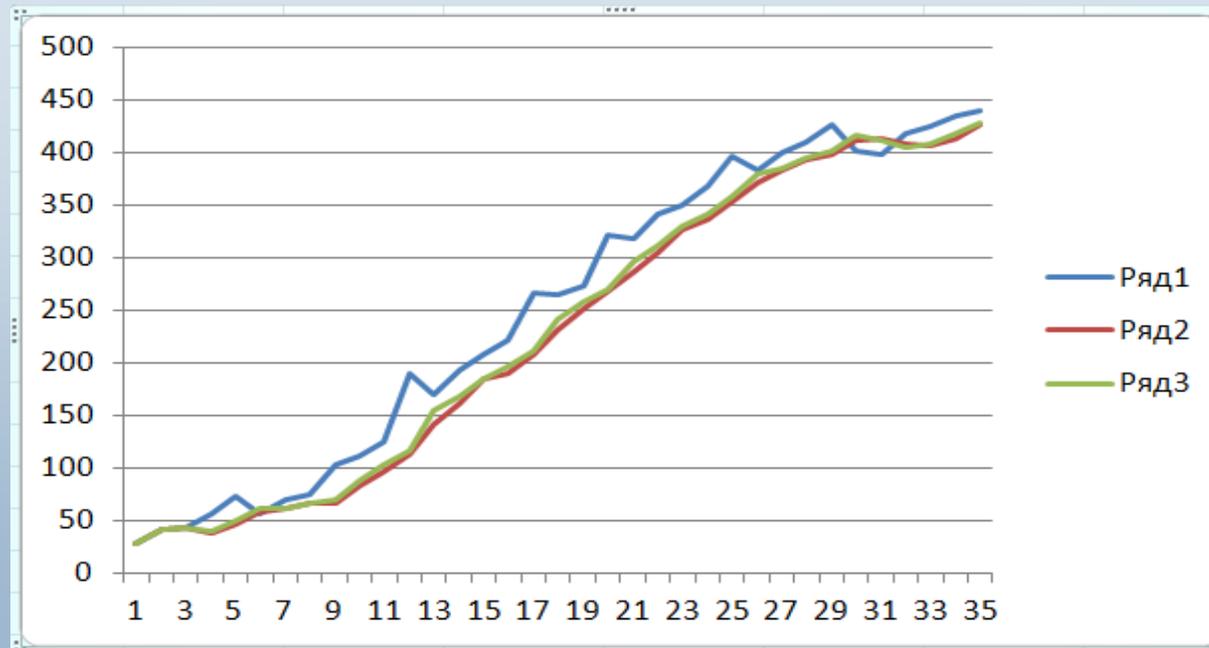
Взвешенное скользящее среднее.

Метод взвешенного скользящего среднего заключается в том что предсказанное значение вычисляется как среднее арифметическое по трем предыдущим данным, при этом каждое из трех величин умножается на веса. Сумма весов должна быть равна единице.

Например

$$\hat{y}_7 = \frac{3}{6}y_6 + \frac{2}{6}y_5 + \frac{1}{6}y_4.$$

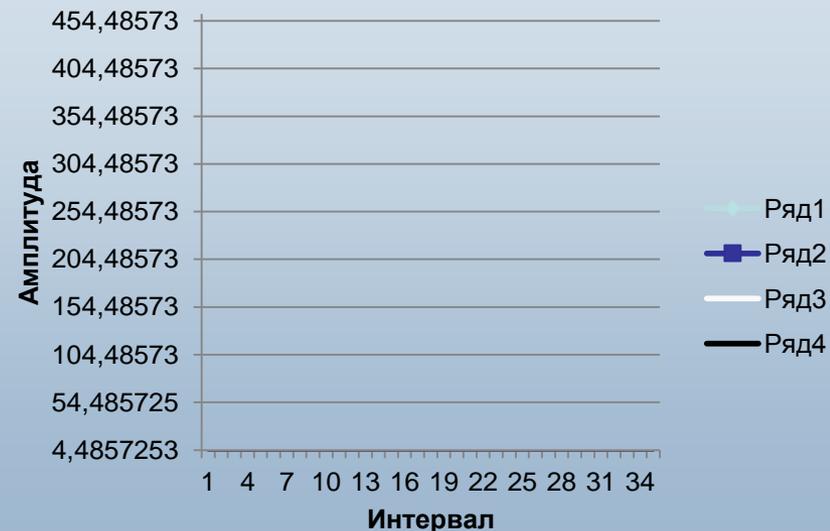
$$\hat{y}_7 = \alpha_1 y_6 + \alpha_2 y_5 + \alpha_3 y_4.$$



Взвешенное скользящее среднее в AtteStat

Увеличение степени полинома улучшает качество аппроксимации

Скользящее среднее



Анализ временных рядов и прогнозирование

Интервал данных: Лист1!\$B\$1:\$B\$35

Интервал вывода: Лист1!\$D\$1

Метод анализа

Скользящее среднее

Сингулярный спектральный анализ

Автокорреляция процесса

Гармонический анализ Фурье

Периодограмма

Сезонный разностный оператор

Параметры сезонного оператора

Период сезонности: 1

Параметры скользящего среднего

Полуширина окна: 1

Степень полинома: 1

Доверительная вероятность: 0,95

Параметры спектрального анализа

Ширина окна: 1

Число гармоник: 1

Параметры гармонического анализа

Число гармоник: 1

Выполнить расчет Отмена Помощь

Экспоненциальное сглаживание (Модель Брауна).

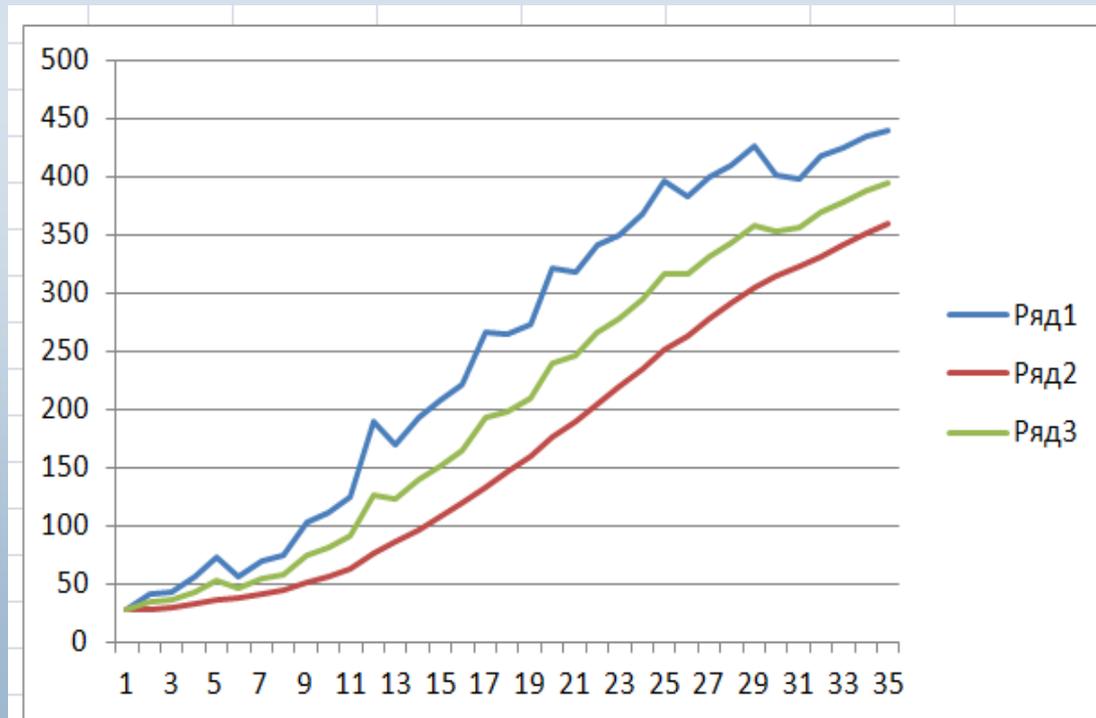
Модель экспоненциального сглаживания отличается от метода взвешенного скользящего среднего выбором коэффициента:

где \hat{Y}_t предсказанное значение за прошлую дату

$$\hat{Y}_{t+1} = \hat{Y}_t + \alpha(Y_t - \hat{Y}_t)$$

Идея метода заключается в том, что прогнозное значение определяется через предыдущее спрогнозированное значение, но скорректированное на величину отклонения факта от прогноза

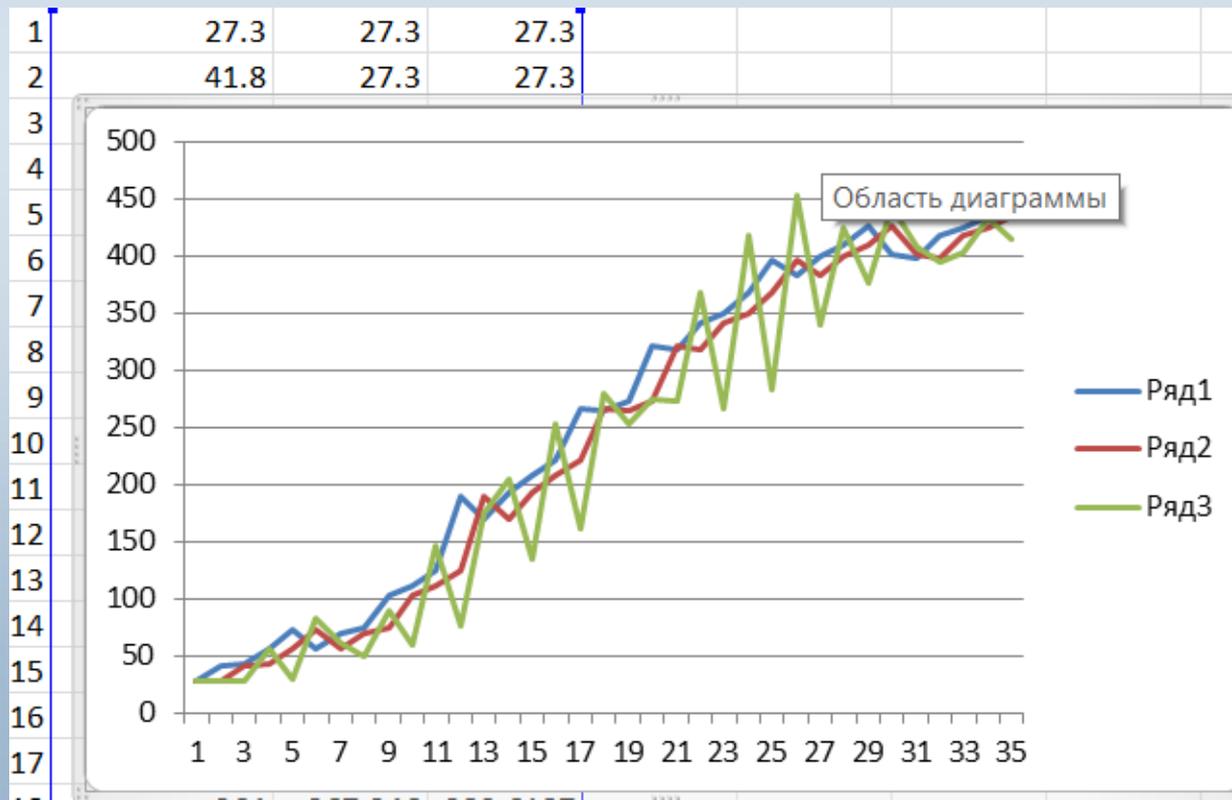
$$\hat{Y}_{t+1} = \alpha Y_t + (1 - \alpha)\hat{Y}_t$$



Экспоненциальное сглаживание (Модель Брауна).

Вообще же модель Брауна может применяться в двух случаях:

1. Когда нужно сгладить имеющийся ряд данных для выявления какой-либо тенденции (обычно в случае со стационарными процессами). Тогда обычно исследователь задаёт значение α в пределах от 0 до 1.

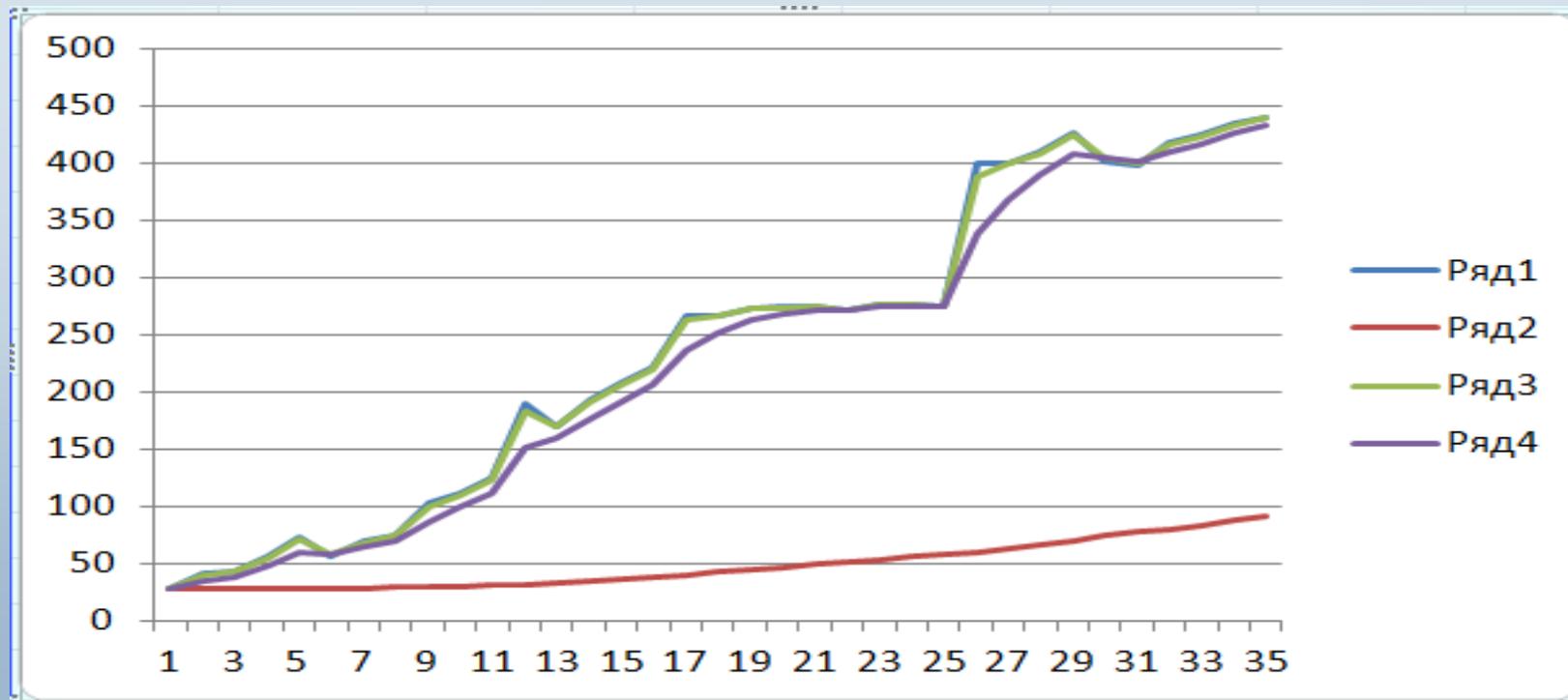


Экспоненциальное сглаживание (Модель Брауна).

Очевидно, что величина α сильно влияет на результат сглаживания. Рассмотрим как метод сглаживания влияет на три ситуации.

1. Реакция сглаживания на скачок.

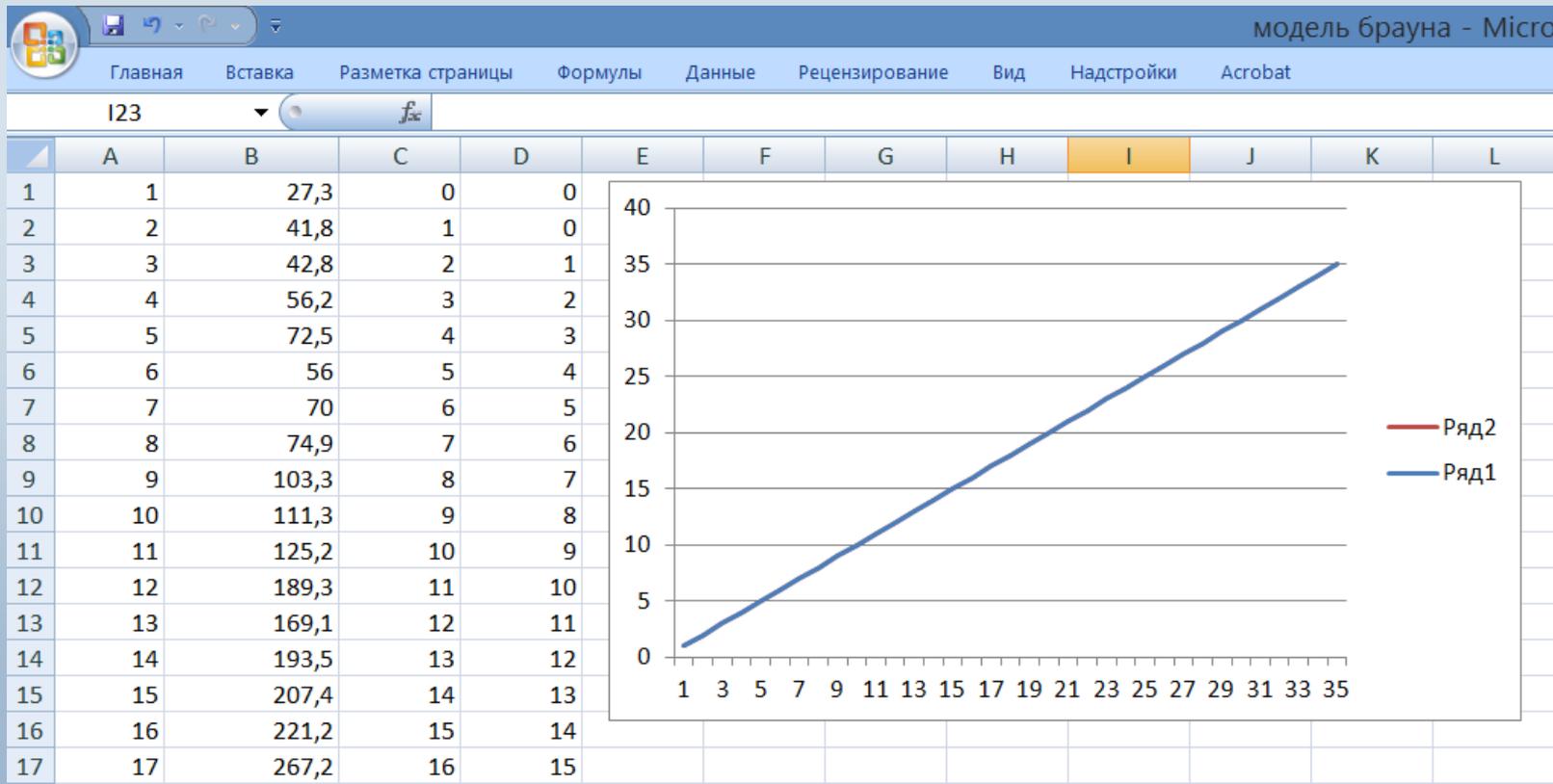
Чем меньше значение α тем хуже результат (ряд 2).



Экспоненциальное сглаживание (Модель Брауна).

1. Реакция сглаживания на постоянные изменения.

Модель Брауна плохо подходит для временных рядов, прогноз всегда занижен.



Выбор адекватной модели прогнозирования

Какую из описанных моделей следует применять для прогнозирования значения временного ряда? Есть несколько эмпирических правил:

- 1. Оценить величину остаточной ошибки с помощью квадратов разностей.** (то есть оценить величину дисперсии). Если модель идеально аппроксимирует значения временного ряда в предыдущие моменты времени, стандартная ошибка оценки равна нулю. С другой стороны, если модель плохо аппроксимирует значения временного ряда в предыдущие моменты времени, стандартная ошибка оценки велика. Таким образом, анализируя адекватность нескольких моделей, можно выбрать модель, имеющую минимальную стандартную ошибку оценки.

Основным недостатком такого подхода является преувеличение ошибок при прогнозировании отдельных значений. Иначе говоря, любая большая разность между величинами Y_i и \hat{Y}_i при вычислении суммы квадратов ошибок SSE возводится в квадрат, т.е. увеличивается.

Выбор адекватной модели прогнозирования

2. Оценить величину остаточной ошибки с помощью абсолютных разностей.

$$MAD = \frac{\sum_{i=1}^n |Y_i - \hat{Y}_i|}{n}$$

При анализе конкретных моделей величина MAD представляет собой среднее значение модулей разностей между фактическим и предсказанными значениями временного ряда. Если модель идеально аппроксимирует значения временного ряда в предыдущие моменты времени, среднее абсолютное отклонение равно нулю. С другой стороны, если модель плохо аппроксимирует такие значения временного ряда, среднее абсолютное отклонение велико. Таким образом, анализируя адекватность нескольких моделей, можно выбрать модель, имеющую минимальное среднее абсолютное отклонение.

Выбор адекватной модели прогнозирования

3. Руководствоваться принципом экономии

Если анализ стандартных ошибок оценок и средних абсолютных отклонений не позволяет определить оптимальную модель, можно воспользоваться четвертым методом, основанным на принципе экономии. Этот принцип утверждает, что из нескольких равноправных моделей следует выбирать простейшую.

