

Множественный регрессионный анализ – это метод математической статистики, который позволяет найти наиболее точное и достоверное отображение (модель, аппроксимацию, уравнение регрессии) стохастической зависимости между откликом  $Y$  и несколькими факторами  $X_1, X_2, \dots, X_j, \dots, X_m$ . Для решения данной задачи необходимо:

1. Определить вид уравнения регрессии.
2. Оценить значимость коэффициентов регрессии.
3. Оценить допустимость отображения исследуемой зависимости выбранным уравнением регрессии.
4. Исследовать остатки (отклонения действительных значений отклика от предсказываемых по уравнению регрессии).

## 8.1. Определение вида уравнения множественной регрессии

Учитывая возможные отклонения, модель связи отклика с некоторым комплексом факторов  $\vec{X} = \{X_1, \dots, X_j, \dots, X_m\}$  должна быть представлена в виде двух составляющих:

$$y = \varphi(\vec{X}) + \varepsilon, \quad (8.1)$$

где  $\varphi(\vec{X})$  – систематическая (объясненная) составляющая. Она обусловлена существованием зависимости между откликом и комплексом факторов;

$\varepsilon$  – случайная составляющая. Она обусловлена разнообразными возмущениями и вызывает отклонения  $y$  от значений, соответствующих реальной зависимости.

Для построения множественной регрессионной модели (иначе – множественного уравнения регрессии или просто множественной регрессии) необходимо решить следующие задачи:

Задача определения вида уравнения множественной регрессии состоит в нахождении систематической составляющей  $\varphi(\vec{X})$ . Однако, поскольку используются выборки ограниченного

объема ( $n \ll \infty$ ), могут быть найдены лишь оценки истинных параметров.

Пусть, например, действительная зависимость отклика от комплекса факторов является линейной:

$$y = \varphi(\vec{X}) = \beta_0 + \sum_{j=1}^m \beta_j X_j. \quad (8.2)$$

Оценкой (моделью, отображением, аппроксимацией) этой связи также может быть линейное выражение:

$$\hat{y} = \hat{\varphi}(\vec{X}) = b_0 + \sum_{j=1}^m b_j X_j. \quad (8.3)$$

В выражении (8.3), которое и есть уравнение регрессии, коэффициенты регрессии  $b_0$  и  $b_j$  ( $j=1, 2, \dots, m$ ) представляют собой оценки коэффициентов истинной зависимости ( $b_0 \approx \beta_0$  и  $b_j \approx \beta_j$ ).

Для подбора уравнения  $\hat{y} = \hat{\varphi}(\vec{X})$ , которое наилучшим образом отображает стохастическую связь между откликом и рассматриваемыми факторами, используют метод наименьших квадратов (МНК). Согласно МНК наилучшей оценкой исследуемой зависимости является та, которая дает наименьшую сумму квадратов отклонений наблюдаемых значений отклика  $y_i$  от рассчитанных по уравнению регрессии  $\hat{y}_i$  при тех же значениях факторов  $\{x_{1i}, \dots, x_{ji}, \dots, x_{mi}\}$ . Это условие выражается следующим образом:

$$\sum_{i=1}^n e_i^2 = \sum_{i=1}^n (y_i - \hat{y}_i)^2 = \min. \quad (8.4)$$

Исходя из условия (8.4) задача определения коэффициентов уравнения регрессии сводится практически к определению минимума функции нескольких переменных и решена математической статистикой для линейного уравнения. Значения коэффициентов регрессии в (3) вычисляются решением системы из  $n$  линейных уравнений с  $m$  неизвестными (здесь  $n$  - число наблюдений) [1, 2, 3].

В *MS Excel* коэффициенты линейной аппроксимации могут быть определены с использованием статистической функции ЛИНЕЙН(). Синтаксис функции и вопросы, связанные с ее использованием приведены в прил.2. Для комплексного решения задачи

множественного регрессионного анализа в *MS Excel* имеется инструмент «Регрессия».

## 8.2.Оценивание качества множественной аппроксимации

Для оценивания качества множественной аппроксимации необходимо определить ее статистическую надежность, проверить значимость коэффициентов регрессии и проверить выполнение допущений относительно остатков. Оценивание статистической надежности уравнения множественной регрессии выполняется так же, как и парной аппроксимации (см. подраздел 6.). Поэтому здесь рассмотрим только две последние задачи.

### 8.2.1.Проверка значимости коэффициентов регрессии

Коэффициенты регрессии  $b_j$  являются случайными величинами с математическими ожиданиями  $\beta_j$  и дисперсиями, которым соответствуют стандартные отклонения  $S_{bj}$ . Значение  $b_j$  признается статистически значимым, если выполняется условие:

$$t_{bj} = \frac{|b_j|}{S_{bj}} > t[\alpha; n - k], \quad (8.5)$$

где  $t_{bj}$  и  $t[\alpha; n - k]$  - расчетное и табличное числа Стьюдента.

Если условие (5) не выполнено, то следует признать, что  $b_j = 0$  и влияние фактора  $X_j$  на отклик несущественное. В таком случае рекомендуется повторить регрессионный анализ без учета фактора  $X_j$ .

### 8.2.2.Анализ остатков

Остатками принято называть отклонения действительных значений отклика от рассчитанных по уравнению регрессии (для  $i$ -го наблюдения остаток  $e_i = y_i - \hat{y}$ ). Отклонения обусловлены наличием в (1) случайной составляющей  $\varepsilon$ , относительно которой делают следующие предположения:

1. Это нормально распределенная случайная переменная.
2. Математическое ожидание случайной составляющей равно нулю -  $M(\varepsilon) = 0$ . Считают, что данная гипотеза выполняется, если среднее выборочное остатков можно считать равным нулю:

$$\bar{e} = \frac{1}{n} \sum_{i=1}^n e_i \approx 0. \quad (8.6)$$

3. Дисперсия случайной составляющей постоянна -  $D(\varepsilon) = Const$ . Гипотеза может быть проверена, например, построением графиков остатков в зависимости от каждого фактора. Если на всех таких графиках остатки примерно равномерно рассеяны в пределах области, параллельной оси  $X_j$  (рис. П1.1а), то гипотезу  $D(\varepsilon) = Const$  считают справедливой.
4. В различных наблюдениях значения  $\varepsilon$  не зависят друг от друга. Для проверки гипотезы о независимости отклонений в различных наблюдениях оценивают автокорреляцию остатков с применением критерия Дарбина-Уотсона (критерия  $DW$ ):

$$DW = \frac{\sum_{i=2}^n (e_i - e_{i-1})^2}{\sum_{i=1}^n e_i^2}. \quad (8.7)$$

Строгое условие отсутствия автокорреляции  $DW = 2$  [2, 6]. Однако, с учетом особенностей распределения критерия Дарбина-Уотсона, ориентировочно можно считать, что автокорреляция остатков отсутствует при  $1,2 \leq DW \leq 2,8$  [6]. В противном случае следует признать, что гипотеза о независимости остатков в рассматриваемом случае не верна.

Если анализ остатков обнаруживает несоответствия указанным гипотезам, то уравнение регрессии, относительно которого данные остатки получены, следует считать неудовлетворительным, т. к. правомерность применения МНК и указанных выше оценок и для множественного регрессионного анализа может быть поставлена под сомнение. В таком случае рекомендуют рассмотреть уравнение иного вида (например, нелинейное вместо линейного), включить неучтенные ранее факторы, выделить в области варьирования факторов различные подобласти.

### **8.3. Пример множественного регрессионного анализа в *MSE Excel* с применением инструмента «Регрессия»**

В качестве примера рассмотрим задачу построения зависимости предела текучести металла, прокатанного на широко-

	A	B	C	D	E	F	G	H	I	J	K	L	M	N
1	<b>ИСХОДНЫЕ ДАННЫЕ</b>				<b>ВЫВОД ИТОГОВ</b>									
2	St, МПа	тыс.С	том.С											
3	394,0	788	560		Регрессионная статистика									
4	341,2	788	584		Множественный R	0,937								
5	338,4	788	611		R-квадрат	0,879								
6	294,1	788	639		Нормированный R-квадрат	0,869								
7	330,9	788	657		Стандартная ошибка	18,56								
8	281,9	788	715		Наблюдения	28								
9	285,9	788	728											
10	336,4	834	550		Дисперсионный анализ									
11	320,5	834	584			df	SS	MS	F	Значимо сть F				
12	336,8	834	611		Регрессия	2	62288,036	31134,018	90,4187	3,6E-12				
13	296,4	834	639		Остаток	25	8608,291	344,332						
14	286,3	834	657		Итого	27	70876,327							
15	240,3	834	715											
16	206,3	834	728			Коэффици	Стандартная	t	P	Набные	Верхние	Набные	Верхние	
17	334,2	875	560		ценный	1152,695	ошибка	статистика	Значение	95%	95%	90,0%	90,0%	
18	290,3	875	584		У-пересечение	73,484		15,8863	1,9E-14	1001,362	1304,039	1027,174	1278,217	
19	318,5	875	611		тыс.С	-0,468		0,073	-6,7941	4,03E-07	-0,649	-0,347	-0,623	-0,373
20	283,2	875	639		том.С	-0,665		0,060	-11,6051	1,47E-11	-0,818	-0,572	-0,797	-0,593
21	281,0	875	657											
22	233,1	875	715											
23	200,0	875	728		<b>ВЫВОД ОСТАТКА</b>									
24	304,4	917	560											
25					Наблюдение	Предска занное Ст.	Остатки	$(e_i - e_{i-1})^2$						
26	307,5	917	584		1	371,3	-7,3	53,6						
27	258,8	917	611		2	354,6	-13,4	37,48						
28	270,3	917	657		3	335,9	2,5	254,80						

Рис. 8.1. Фрагмент рабочего листа с результатами множественного регрессионного анализа инструментом «РЕГРЕССИЯ»

полосном стане горячей прокатки (ШСГП) от температур конца прокатки (ткп) и смотки (тсм). Исходные данные заносятся на рабочий лист с клавиатуры (на рис. 8.1 и 8.2 они расположены в ячейках A1:C29).

### 8.3.1. Аппроксимация с применением инструмента

#### «РЕГРЕССИЯ»

При выполнении работы настройки инструмента «Регрессия» должны соответствовать указанным в приложении 1.

В этом случае основные результаты (ячейки E1:M19 на рис. 8.1) будут дополнены таблицей остатков (ячейки E23:G23 на рис. 8.2).

	A	B	C	D	E	F	G	H
22	233,1	875	715					
23	200,0	875	728		ВЫВОД ОСТАТКА			
24	304,4	917	560					
25	307,5	917	584		<i>Наблюдение</i>	<i>Предсказанное Значение Ст, МПа</i>	<i>Остатки</i>	$(e_i - e_{-1})^2$
26	268,8	917	611		1	371,3	-7,3	53,6
27	225,4	917	639		2	354,6	-13,4	37,48
28	270,3	917	657		3	335,9	2,5	254,80
29	181,3	917	715		4	316,4	-32,3	1214,00
30	175,3	917	728		5	303,9	27,0	3517,48
31					6	263,6	18,3	75,61
32					7	254,6	11,3	48,53
33					8	348,4	-12,0	545,15
34					9	331,8	-11,3	0,60
35					10	313,0	23,8	1229,38
36					11	293,5	5,9	321,94
37					12	281,0	-11,7	309,47
38					13	240,7	-0,4	127,80
39					14	231,7	-25,4	623,31
40					15	328,0	6,2	996,21
41					16	311,3	-21,0	741,05
42					17	292,6	25,9	2205,48
43					18	273,1	10,1	250,99
44					19	260,6	20,4	106,26
45					20	220,3	12,8	57,69
46					21	211,3	-11,3	579,18
47					22	307,1	-2,7	73,28
48					23	290,4	17,1	391,16
49					24	271,7	-2,9	397,50
50					25	252,2	-26,8	573,25
51					26	239,7	30,6	3295,72
52					27	199,4	-18,1	2371,23
53					28	190,4	-15,1	9,20
54								

Рис. 8.2. Таблица остатков, полученная с использованием инструмента «РЕГРЕССИЯ»

Таблицу остатков необходимо дополнить столбцом, во всех строках которого, начиная со второй, вычислить значения  $(e_i - e_{i-1})^2$ . Эти данные будут в дальнейшем использованы для расчета критерия  $DW$ . Например, в ячейке H27:

$$=(G27-G26)^2.$$

На основании результатов работы инструмента записать уравнение регрессии в содержательной форме, оценить значимость коэффициентов регрессии, надежность аппроксимации и автокорреляцию остатков.

**Уравнение регрессии в содержательной форме** (ячейки H1:J2) записывается с клавиатуры.

**Оценивание значимости коэффициентов регрессии** (ячейки H3:J6). В ячейке H5 определяется табличное число Стьюдента. Для рассматриваемого примера (число коэффициентов регрессии  $k=3$ ):

$$=\text{СТЮДРАСПОБР}(0,05;2;F8-3).$$

В ячейках J4:J6, с использованием функции ЕСЛИ() программируется вывод о значимости коэффициентов регрессии. Например, для ячейки J4:

$$=\text{ЕСЛИ}(\text{ABS}(H17)>\$H\$5;"Значим";"Не значим").$$

**Внимание!** Если коэффициент регрессии при факторе  $X_j$  оказался не значимым, необходимо повторить регрессионный анализ, не включая во входной интервал  $X$  столбец со значениями  $X_j$ . Поскольку входной интервал  $X$  должен состоять из смежных столбцов, может оказаться необходимым перегруппировать столбцы со значениями факторов.

**Оценивание надежности аппроксимации** в примере на рис. 8.1 выполняется в ячейках H7:J9. В ячейке H9, с помощью статистической функции ФРАСПОБР(), определяется табличное значение числа Фишера:

$$=\text{ФРАСПОБР}(0,05;2;F8-3).$$

Слово «Аппроксимация» в ячейки I8:J8 введено с клавиатуры. Собственно вывод («надежная» или «не надежная») формируется в ячейке I9 с применением функции ЕСЛИ() путем сравнения табличного (из ячейки H9) и рассчитанного (из ячейки I12) чисел Фишера:

$$=\text{ЕСЛИ}(H9>I12;"надежная";"не надежная").$$

**Оценивание автокорреляции остатков** выполнено в ячейках L7:M9. Значение критерия Дарбина-Уотсона в ячейке M8 вычисляется по формуле (10), которая для рассматриваемого примера программируется следующим образом:

$$= \text{СУММ}(H27:H53)/G13.$$

Вывод в ячейке L9 формируется с применением функции ЕСЛИ(), которая проверяет условие  $1,2 \leq DW \leq 2,8$ :

ЕЛИ(M8<1,2;"Существует";ЕСЛИ(M8>2,8;"Существует";"Отсутствует")) .

#### 8.4. Анализ результатов множественного регрессионного анализа

Анализируя результаты множественного регрессионного анализа необходимо ответить на следующие вопросы:

1. Связь между какими величинами анализировалась?
2. Значимы ли коэффициенты регрессии?
3. Как выглядит уравнение множественной регрессии, полученное в результате выполнения работы?
4. Можно ли считать полученное уравнение множественной регрессии статистически надежной аппроксимацией анализируемой зависимости?

Применительно к рассмотренному примеру можно сказать следующее. Анализировалась связь между пределом текучести металла  $\sigma_T$ , температурой конца прокатки  $t_{кп}$  и смотки  $t_{см}$  на ШСГП.

С доверительной вероятностью  $p = 95\%$  коэффициенты регрессии  $b(t_{кп}) = -0,498$  и  $b(t_{см}) = -0,695$  являются статистически значимыми, т. к. соответствующие числа Стьюдента  $|t(t_{кп})| = 6,7941$  и  $|t(t_{см})| = 11,6051$  больше табличного  $t[0,05;25] = 2,0595$ .

Для рассмотренных условий множественная линейная аппроксимация зависимости предела текучести металла от температуры конца прокатки и смотки на ШСГП имеет вид:

$$\sigma_T = 1152,695 - 0,498 t_{кп} - 0,695 t_{см}.$$

С доверительной вероятностью 95% полученное уравнение регрессии можно считать статистически надежной аппроксимацией исследуемой зависимости, т. к. рассчитанное число Фишера  $F_p = 90,4187$  больше табличного  $F[0,05;2;25] = 3,3352$ .

#### 8.5. Контрольные вопросы

1. Поясните сущность и укажите этапы множественного регрессионного анализа.
2. Укажите допущения множественного регрессионного анализа.
3. Запишите модель множественного регрессионного анализа.
4. Что представляет собой уравнение множественной регрессии?
5. Как определить качество уравнения множественной регрессии?