

Парный регрессионный анализ – это метод математической статистики, который позволяет найти наиболее точное и достоверное отображение (модель, аппроксимацию, уравнение регрессии) стохастической зависимости между откликом Y и одним из факторов X . Для решения данной задачи необходимо:

1. Определить вид уравнения регрессии
2. Оценить допустимость отображения исследуемой зависимости выбранным уравнением регрессии

7.1. Содержание и допущения регрессионного анализа

Понятие стохастической зависимости в некотором смысле является обобщением понятия о зависимости функциональной. Для последней характерно, что каждому значению фактора (аргумента) соответствует совершенно определенное значение отклика. В случае стохастической зависимости при определенном значении x_i фактора X может наблюдаться множество значений отклика Y . В производственных условиях фактор также является случайной величиной, но при проведении регрессионного анализа полагают, что всякое его значение x_i неслучайно.

При проведении регрессионного анализа принимают следующие допущения.

1. Фактор измеряется с пренебрежимо малой ошибкой по сравнению с ошибкой определения отклика. Большая ошибка y_i объясняется наличием в каждом из наблюдений влияний на отклик нерегламентированных параметров.

2. Каждому значению фактора соответствуют значения отклика, представляющие собой независимые и нормально распределенные величины.

3. При проведении эксперимента с объемом выборки n при условии, что каждый опыт произведен m раз, выборочные дисперсии отклика однородны.

Учитывая возможные отклонения, модель связи некоторого значения отклика с соответствующим значением фактора может быть представлена в виде двух составляющих:

$$y_i = \varphi(x_i) + \varepsilon_i, \quad (7.1)$$

где $\varphi(x_i)$ - систематическая (объясненная) составляющая. Она обусловлена существованием зависимости между откликом и фактором;

ε_i - случайная составляющая. Она обусловлена разнообразными возмущениями и вызывает отклонения y_i от соответствующих реальной зависимости.

Относительно ε_i делают следующие предположения:

1. Это нормально распределенная случайная переменная.
2. $\mu(\varepsilon_i) = 0$ (математическое ожидание случайной составляющей равно нулю).
3. $\sigma(\varepsilon_i) = Const$ (дисперсия случайной составляющей постоянна).
4. В различных наблюдениях значения ε_i не зависят друг от друга.

Для построения парной регрессионной модели (иначе – парного уравнения регрессии или просто парной регрессии) необходимо решить две задачи - определить вид уравнения регрессии и оценить допустимость отображения исследуемой зависимости выбранным уравнением регрессии.

7.2. Определение вида уравнения парной регрессии

Задача определения вида уравнения регрессии состоит в определении систематической составляющей $\varphi(x)$. Однако, как уже указывалось ранее, истинные параметры (коэффициенты) этого уравнения не могут быть определены, поскольку используются выборки ограниченного объема ($n \ll \infty$). Поэтому могут быть найдены лишь оценки истинных параметров и действительная связь между откликом и фактором $y = \varphi(x)$ представляется оценкой (отображением) этой связи $\hat{y} = \hat{\varphi}(x)$. Именно данная оценка и является уравнением регрессии.

Для подбора уравнения $\hat{y} = \hat{\varphi}(x)$, которое наилучшим образом отображает стохастическую связь между откликом и фактором, 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 фактора x_i (рис. 7.1). Это условие выражается следующим образом:

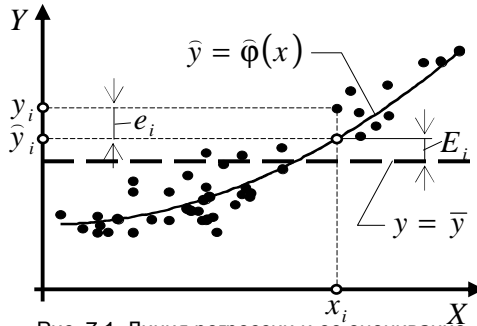


Рис. 7.1. Линия регрессии и ее оценивание

$$\sum_{i=1}^n e_i^2 = \sum_{i=1}^n (y_i - \hat{y}_i)^2 = \min. \quad (7.2)$$

Оценки, полученные МНК, являются несмещенными, состоятельными и эффективными. Несмещенность означает, что математическое ожидание каждого из оцениваемых параметров равно соответствующему истинному значению. Состоятельность означает, что с увеличением числа наблюдений n оценки параметров все более концентрируются вокруг истинных значений (т. е. с возрастанием n дисперсии оценок стремятся к нулю). Эффективность означает, что оценки, полученные МНК, обладают наименьшей дисперсией по сравнению с оценками этих же параметров, полученными другими методами.

Задача определения коэффициентов уравнения регрессии сводится практически к определению минимума функции нескольких переменных и решена математической статистикой для линейного уравнения

$$y = \beta_0 + \beta_1 x \approx b_0 + b_1 x. \quad (7.3)$$

Оценки коэффициентов $b_0 \approx \beta_0$ и $b_1 \approx \beta_1$ вычисляются по формулам (7.4, 7.5):

$$b_1 = \frac{\sum_{i=1}^n x_i \sum_{i=1}^n y_i - n \sum_{i=1}^n x_i y_i}{\left(\sum_{i=1}^n x_i \right)^2 - n \sum_{i=1}^n x_i^2}; \quad (7.4)$$

$$b_0 = \frac{1}{n} \left(\sum_{i=1}^n y_i - b_1 \sum_{i=1}^n x_i \right). \quad (7.5)$$

В MS Excel коэффициенты b_0 и b_1 линейной регрессии можно определить с использованием статистических функций НАКЛОН(<Y>;<X>) и ОТРЕЗОК(<Y>;<X>). Аргументы <Y> и <X> - ссылки на ячейки, где записаны значения отклика и фактора соответственно.

При необходимости вычисления коэффициентов уравнений регрессии других видов поступают следующим образом. Сначала по экспериментальным данным x_i и y_i находят значения x'_i и y'_i , которые соответствуют уравнению данного вида после его линеаризации. Далее используют формулы (7.4) и (7.5), расчеты по которым дадут коэффициенты именно линеаризованного уравнения (b'_0 и b'_1). После этого находят коэффициенты реальных уравнений. Формулы для таких преобразований представлены в табл. 7.1.

Таблица 7.1

Формулы для вычисления коэффициентов некоторых нелинейных уравнений парной регрессии

Аппроксимация	Преобразование исходных данных		Преобразование коэффициентов	
	x'_i	y'_i	b_0	b_1
$b_0 + b_1 / x$	$1 / x_i$	y_i	b'_0	b'_1
$b_0 + b_1 x^m$	x_i^m	y_i	b'_0	b'_1
$b_0 + b_1 \lg x$	$\lg x_i$	y_i	b'_0	b'_1
$b_0 + b_1 \ln x$	$\ln x_i$	y_i	b'_0	b'_1
$1 / (b_0 + b_1 x)$	x_i	$1 / y_i$	b'_0	b'_1
$x / (b_0 + b_1 x)$	x_i	x_i / y_i	b'_0	b'_1
$b_0 / (b_1 + x)$	x_i	$1 / y_i$	$1 / b'_0$	b'_0 / b'_1
$b_0 x / (b_1 + x)$	$1 / x_i$	$1 / y_i$	$1 / b'_0$	b'_0 / b'_1
$b_0 b_1^x$	x_i	$\lg y_i$	$10 b'_0$	$10 b'_1$
$b_0 x^{b_1}$	$\ln x_i$	$\ln y_i$	$\exp(b'_0)$	b'_1
$b_0 \exp(b_1 x)$	x_i	$\ln y_i$	$\exp(b'_0)$	b'_1
$b_0 \exp(b_1 / x)$	$1 / x_i$	$\ln y_i$	$\exp(b'_0)$	b'_1

7.3. Оценивание уравнения парной регрессии

Из различных уравнений регрессии наилучшим в смысле МНК считают то, которое обеспечивает минимум дисперсии фактических (полученных экспериментально) значений отклика относительно линии регрессии. Эту дисперсию называют остаточной или дисперсией относительно регрессии и определяют по формуле:

$$S_e^2 = \frac{1}{n-k} \sum_{i=1}^n (y_i - \hat{y}_i)^2, \quad (7.6)$$

где k - число коэффициентов регрессии в уравнении.

Точность отображения (аппроксимации) исследуемой зависимости выбранным уравнением регрессии оценивают с помощью дисперсионного анализа. Для этого сравнивают дисперсию относительно линии регрессии (S_e^2) с оценкой дисперсии значений y_i

относительно выборочного среднего фактических значений отклика \bar{y} :

$$S_E^2 = \frac{1}{k-1} \sum_{i=1}^n (y_i - \bar{y})^2 = \frac{1}{m} \sum_{i=1}^n \left(y_i - \frac{1}{n} \sum_{i=1}^n y_i \right)^2. \quad (7.7)$$

Величина S_E^2 характеризует рассеяние y_i , обусловленное зависимостью отклика от факторов, и поэтому называется объясненной дисперсией. Остаточная дисперсия S_e^2 характеризует рассеяние y_i , вызванное случайными воздействиями (возмущениями). Очевидно, что связь между откликом и факторами в виде данного уравнения регрессии существует, если объясненная дисперсия существенно больше остаточной.

Чтобы выяснить, можно ли считать отличие рассматриваемых дисперсий существенным, выдвигают нулевую гипотезу об их равенстве $H_0 : S_E^2 = S_e^2$ и проверяют ее с использованием числа Фишера:

$$F = \frac{S_E^2}{S_e^2}. \quad (7.8)$$

Гипотеза считается справедливой, если рассчитанное число Фишера не превышает табличного значения $F[\alpha, \nu_1, \nu_2]$ для заданного уровня значимости α . При выборе табличного значения важно помнить, что в данном случае $\nu_1 = \nu_E = k - 1$ и $\nu_2 = \nu_e = n - k$. В *MS Excel* табличное число Фишера может быть найдено с применением статистической функции:

ФРАСПОБР($\alpha; k - 1; n - k$).

Справедливость гипотезы о равенстве остаточной и объясненной дисперсий означает, что выбранное уравнение регрессии нельзя принимать в качестве модели связи между откликом и фактором.

Если же $F > F[\alpha, k - 1, n - k]$, то объясненная дисперсия существенно больше остаточной. А это означает, что между откликом и факторами существует взаимосвязь, которую с вероятностью α допустимо аппроксимировать рассматриваемым уравнением регрессии.

Из нескольких допустимых аппроксимаций наиболее точной, очевидно, будет та, для которой значение S_e^2 является наименьшим. Отсюда следует, что для наиболее точной модели $\hat{y}_i = \hat{\Phi}(x_i)$ различие расчетного и табличного чисел Фишера будет максимальным.

Во многих программных средствах, содержащих опции обработки данных (в том числе и в *MS Excel*), для оценивания качества аппроксимации предлагается параметр R^2 (коэффициент достоверности аппроксимации):

$$R^2 = 1 - \frac{\sum_{i=1}^n e_i^2}{\sum_{i=1}^n (y_i - \bar{y})^2} \frac{n-1}{n-k}. \quad (7.9)$$

Коэффициент достоверности представляет собой долю дисперсии отклика, объясненную с помощью оцениваемого уравнения регрессии. Чем меньше остаточная дисперсия по отношению к дисперсии значений отклика относительно его среднего выборочного, тем точнее модель $\hat{y} = \hat{\Phi}(x)$ отображает изменчивость Y в связи с изменениями фактора

X , тем больше значение R^2 . Идеальная модель ($\sum_{i=1}^n e_i = 0$) обеспечивает $R^2 = 1$ и в этом случае необходимость в оценивании адекватности аппроксимации отпадает. В остальных случаях следует проверить условие отличия коэффициента детерминации от нуля:

$$F = \frac{R^2}{1 - R^2} \frac{n - k}{k - 1} > F[\alpha; k - 1; n - k]. \quad (7.10)$$

Если при заданной доверительной вероятности указанное условие выполняется, можно считать, что R^2 значимо отличается от нуля и, сле-

довательно, оцениваемое уравнение регрессии является достаточно точной аппроксимацией исследуемой зависимости. В противном случае необходимо признать, что аппроксимация не является адекватной и рассмотреть иной вариант уравнения регрессии.

7.4. Парный регрессионный анализ в MS Excel с применением инструмента «Линия тренда»

7.4.1. Выполнение парного регрессионного анализа с применением инструмента «Линия тренда»

В качестве примера рассмотрим построение аппроксимации зависимости предела текучести металла, прокатанного на ШСП, от температуры смотки тсм. На рабочем листе (рис. 7.2) исходные данные расположены в ячейках B2:B30 и D2:D30.

Сначала с использованием инструмента «Мастер диаграмм», исходные данные необходимо отобразить в виде диаграммы типа «Точечная». Диаграмма должна быть отформатирована таким образом, чтобы ее оформление соответствовало правилам графического представления данных.

После форматирования диаграммы, с помощью мыши выделить на диаграмме ряд данных. Через пункт <Диаграмма> в строке главного меню MS Excel, задать команду <Добавить линию тренда>. На экране появится диалоговое окно <Линия тренда>.

На закладке «Параметры» включить опции «Показывать уравнение на диаграмме» и «Поместить на диаграмму величину достоверности аппроксимации (R^2)». Затем на закладке «Тип» выбрать тип аппроксимации «Линейная». После нажатия кнопки <ОК> на диаграмме появятся линия тренда и надпись, содержащая уравнение регрессии и значение R^2 .

Автоматически фон надписи устанавливается прозрачным, что затрудняет чтение текста. Поэтому рекомендуется придать ей атрибут «С тенью» через закладку «Вид» окна <Формат подписи данных>, которое появиться, если выполнять форматирование надписи через пункт <Формат> главного меню.

Под диаграммой необходимо подготовить обозначения столбцов, в которых будет размещаться информация о рассмотренных уравнениях регрессии. На рис. 7.2 такие обозначения записаны в ячейках I22:O22.

В ячейке J22 с применением функции СЧЕТ() определяется число наблюдений:

$$=СЧЁТ(B3:B30).$$

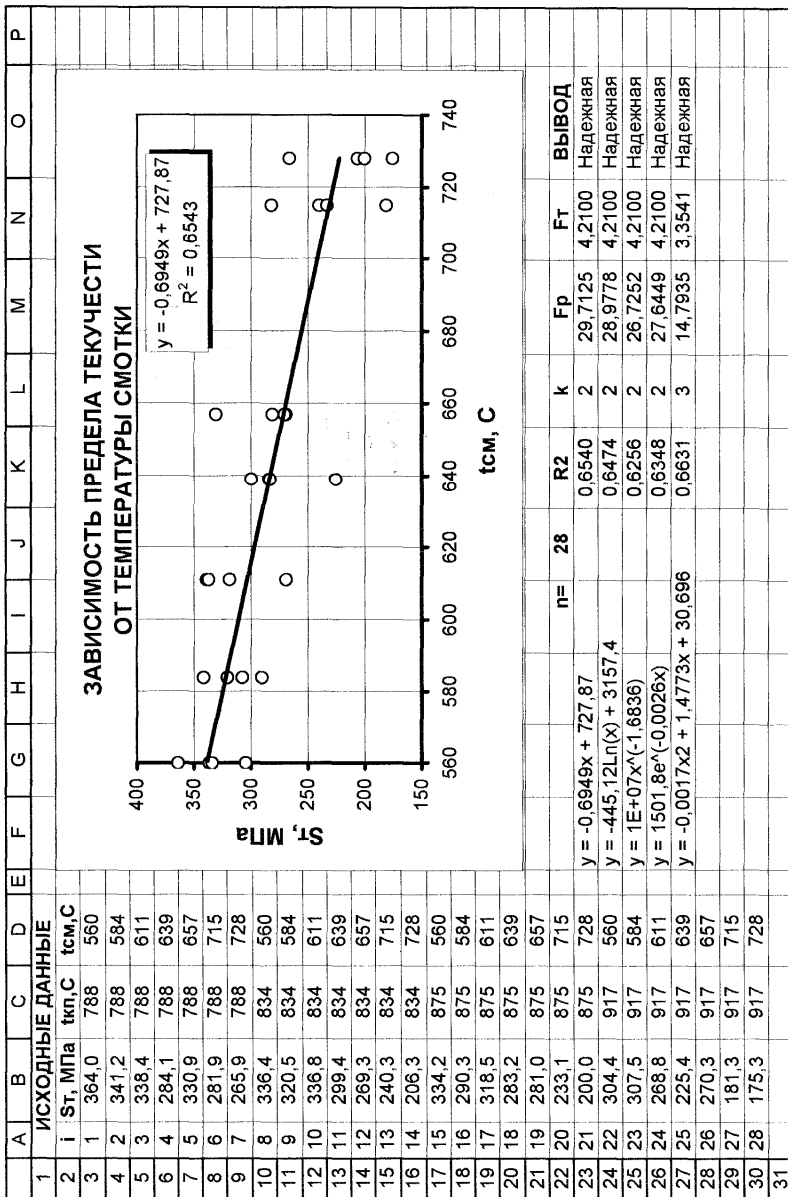


Рис. 7.-2. Фрагмент рабочего листа с исходными данными и результатами парного регрессионного анализа зависимости предела текучести от температуры смотки

Уравнение регрессии и значение R^2 записываются с клавиатуры на основании надписи, которой сопровождается линия тренда. Число коэффициентов регрессии k также вводится с клавиатуры.

Рассчитанное число Фишера F_p вычисляется по формуле (6.10). Например, в ячейке M23:

$$=K23*(\$J\$22-L23)/((1-K23)*(L23-1)).$$

Табличное число Фишера F_{95} находится для доверительной вероятности 95% с использованием статистической функции ФРАСПОБР(). Например, в ячейке N23:

$$=ФРАСПОБР(0,05;L23-1;J\$22-1).$$

Чтобы рассмотреть другие аппроксимации, доступные через инструмент «Линия тренда», необходимо выполнять форматирование линии тренда.

Форматирование может быть реализовано либо через контекстное меню, либо через пункт <Формат> главного меню. В этом случае на экране также появляется окно <Линия тренда>, но теперь достаточно указать другой тип аппроксимации на закладке «Тип». В результате на диаграмме вместо предыдущего тренда появиться новый, а его параметры отобразятся в надписи. **Внимание!** Если каждый из вариантов аппроксимации наносить на диаграмму как добавление линии тренда, то на диаграмме будут одновременно отображены все рассмотренные варианты.

7.4.2. Оценивание результатов анализа зависимости предела текучести от температуры смотки

В выводах по результатам парного регрессионного анализа необходимо ответить на следующие вопросы:

1. Связь между какими величинами анализировалась?
2. Сколько и какие аппроксимации были рассмотрены, какие из них являются статистически значимыми?
3. Какая из аппроксимаций является наилучшим отображением связи между анализируемыми параметрами и почему можно это утверждать?

Относительно результатов рассмотренного примера можем сделать следующие выводы.

1. Анализировалась связь между пределом текучести металла σ_T и температурой смотки t_{cm} при прокатке на ШСГП.
2. Рассмотрели пять аппроксимаций:

$$\begin{aligned} \sigma_T &= 0,6949 t_{cm} + 727,87 & (1) & \quad F_p = 49,1445 \quad F_{95} = 4,2100 \\ \sigma_T &= 3157,4 - 445,12 \ln(t_{cm}) & (2) & \quad F_p = 47,7379 \quad F_{95} = 4,2100 \\ \sigma_T &= 10^7 t_{cm}^{-1,6836} & (3) & \quad F_p = 43,4444 \quad F_{95} = 4,2100 \\ \sigma_T &= 1501,8 e^{-0,0026 t_{cm}} & (4) & \quad F_p = 45,1939 \quad F_{95} = 4,2100 \\ \sigma_T &= 30,696 + 1,4773 t_{cm} - 0,0017 t_{cm}^2 & (5) & \quad F_p = 24,6030 \quad F_{95} = 3,3541 \end{aligned}$$

3. С доверительной вероятностью 95% статистически значимыми являются все рассмотренные аппроксимации, т. к. во всех случаях рассчитанные числа Фишера больше табличных.

4. Наилучшим отображением связи между пределом текучести металла и температурой смотки является линейная аппроксимация

$$\begin{aligned} \sigma_T &= 727,87 - 0,6949 t_{cm} \\ (R^2 &= 0,6540; F_p = 49,1445; F_{95} = 4,2100) \end{aligned}$$

так как для нее характерно наибольшее различие между расчетным и табличным числами Фишера.

7.5. Построение аппроксимации с применением статистических функций

Порядок решения задачи иллюстрируется на примере построения зависимости предела текучести стали от степени деформации при холодном пластическом деформировании.

В практике ОМД для отображения зависимости предела текучести от степени деформации при холодном пластическом деформировании наибольшее распространение получили выражения следующего вида [6, 7]:

$$\sigma_T = \sigma_{T0} + b_0 \varepsilon^{b_1}; \quad (7.11)$$

$$\sigma_T = \sigma_{T0} + b_1 \sqrt{\varepsilon} + b_0, \quad (7.12)$$

где σ_{T0} - предел текучести металла в отожженном (недеформированном) состоянии.

Оба выражения можно представить в иной форме:

$$\sigma_T = \sigma_{T0} + \Delta\sigma_T \quad (7.13)$$

Тогда задача построения зависимости для расчета предела текучести сводится к задаче поиска выражения, отображающего влияние степени деформации ε на приращение предела текучести $\Delta\sigma_T$.

Строится аппроксимация вида $\Delta\sigma_T = b_0 + b_1 \sqrt{\varepsilon}$. Пример оформления рабочего листа представлен на рис. 7.3. Исходные данные расположены в ячейках A5:B21.

Сначала необходимо выполнить линеаризацию исходных данных по формулам, приведенным в приложении 7,9. На рис. 7.3 результаты линеаризации расположены в ячейках N4:O20. Они получены путем преобразований $y'_i = y_i$ и $x'_i = x_i^m$ для $m=0,5$. Например, в ячейках N4 и O4 соответственно запрограммировано:

$$=B5 \text{ и } =A5^{\wedge}Q\$3.$$

Затем следует определить коэффициенты линейной регрессии $y'_i = b'_0 + b'_1 x'_i$. В примере (рис. 7.3) значения b'_0 и b'_1 представлены в ячейках P5 и Q5 соответственно. Для их вычисления здесь использованы статистические функции:

$$=ОТРЕЗОК(N4:N20;O4:O20) \text{ и } =НАКЛОН(N4:N20;O4:O20).$$

Далее необходимо осуществить переход к действительным коэффициентам регрессии. В соответствии с приложением 10 для рассматриваемой аппроксимации переход осуществляется с помощью преобразований $b_0 = b'_0$ и $b_1 = b'_1$. В примере действительные значения коэффициентов регрессии представлены в ячейках P7 и Q7. Здесь соответственно запрограммировано:

$$=P5 \text{ и } =Q5.$$

Таким образом, получена следующая аппроксимация:

$$\Delta\sigma_{\tau} = -109,9 + 57,74\sqrt{\epsilon} \quad (7.14)$$

Теперь с использованием условия (7.10) необходимо выполнить оценку статистической надежности данного уравнения регрессии. Оценивание производится в ячейках P9:O13.

В ячейке O9 определяется число наблюдений с помощью функции СЧЕТ():

$$=СЧЁТ(B5:B21).$$

В ячейке O10 вводится число коэффициентов регрессии. Оно принято равным 3, так как аппроксимация в виде (7.14) на самом деле является иной записью уравнения $\Delta\sigma_{\tau} = b_0 + b_1 \epsilon^m$, где третий коэффициент $m=0,5$.

В ячейках R4:R20 вычисляются остаточные отклонения $e_i = y_i - (b_0 + b_1 x_i^m)$. Например, в ячейке R4 записано:

$$=B5-(P7+Q7*(A5)^Q3).$$

В ячейках S4:S20 вычисляются отклонения $E_i = y_i - \bar{y}$. Например, в ячейке S4 записано:

$$=B5-CPЗНАЧ(B5:B21).$$

В ячейках R21 и S21 рассчитывается суммы квадратов отклонений

$$\sum_{i=1}^n e_i^2 \text{ и } \sum_{i=1}^n E_i^2. \text{ В R21 записано}$$

$$=\text{СУММКВ(R4:R20)},$$

а в S21 записано

$$=\text{СУММКВ(S4:S20)}.$$

Показатель достоверности аппроксимации R^2 рассчитан в ячейке Q11 по формуле (7.9):

$$=1-R21*(Q9-1)/(S21*(Q9-Q10)).$$

В ячейке Q12 вычисляется расчетное число Фишера F_p

$$=Q11*(Q9-Q10)/((1-Q11)*(Q10-1)).$$

Чтобы обеспечить возможность оценивания аппроксимации при различных доверительных вероятностях, предусмотрено задание величины p с клавиатуры в ячейку Q13. Далее это значение воспроизводится в обозначении табличного числа Фишера, которое записывается в ячейке R13 с применением следующей программной конструкции:

$$="F"&\text{ТЕКСТ}(Q13;"00").$$

Собственно значение табличного числа Фишера определяется в ячейке Q13 следующим образом:

$$=\text{ФРАСПОБР}(1-Q13/100;Q10-1;Q9-Q10).$$

Вывод о надежности аппроксимации воспроизводится в ячейках P15:Q16. В ячейке P15 содержится неизменная запись «Аппроксимация». В ячейке P16 с использованием функции ЕСЛИ() запрограммировано

$$=\text{ЕСЛИ}(Q12>Q14;"надежная";"ненадежная").$$

Чтобы отобразить найденную аппроксимацию графически, необходимо использовать диаграмму типа «Точечная». Сначала на нее следует нанести ряд анализируемых данных в виде маркеров без линии. Затем на рабочем листе необходимо создать массив данных, где значениями отклика будут результаты расчета по полученному уравнению регрессии. После этого данный массив отображается на той же диаграмме как самостоятельный ряд в виде сглаженной линии без маркеров. Чтобы обеспечить необходимую плавность отображаемой линии регрессии, расчетный массив должен быть содержать достаточно много значений (не менее 90-100).

В рассматриваемом примере (рис. 7.3) такой массив представлен в ячейках U4:V47. Значения фактора (обжатия) в ячейках U4:U47 заданы в пределах от 4 до 90 с шагом 2, что соответствует пределам на оси абсцисс построенной диаграммы. Оценки отклика в ячейках V4:V47 получены расчетами по анализируемому уравнению регрессии. Например, в ячейке V4 записано:

$$=PQ + Q(U)^Q.$$

Ряд, отображающий линию регрессии задан на диаграмме следующим образом:

Значения X: U4:U47

Значения Y: V4:V47

При форматировании ряда на закладке <Вид> окна <Формат ряда данных> указано «Маркер отсутствует» и «Сглаженная линия».

7.6. Контрольные вопросы

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