

12.3: Linear Regression

Linear Regression

Suppose that a pair $\{X, Y\}$ of random variables has a joint distribution. A value $X(\omega)$ is observed. It is desired to estimate the corresponding value $Y(\omega)$. Obvious there is no rule for determining $Y(\omega)$ unless Y is a function of X . The best that can be hoped for is some estimate based on an average of the errors, or on the average of some function of the errors.

Suppose $X(\omega)$ is observed, and by some rule an estimate $\hat{Y}(\omega)$ is returned. The error of the estimate is $Y(\omega) - \hat{Y}(\omega)$. The most common measure of error is the mean of the square of the error

$$E[(Y - \hat{Y})^2]$$

The choice of the mean square has two important properties: it treats positive and negative errors alike, and it weights large errors more heavily than smaller ones. In general, we seek a rule (function) r such that the estimate $\hat{Y}(\omega)$ is $r(X(\omega))$. That is, we seek a function r such that

$$E[(Y - r(X))^2] \text{ is a minimum.}$$

The problem of determining such a function is known as the *regression problem*. In the unit on [Regression](#), we show that this problem is solved by the conditional expectation of Y , given X . At this point, we seek an important partial solution.

The regression line of Y on X

We seek the best straight line function for minimizing the mean squared error. That is, we seek a function r of the form $u = r(t) = at + b$. The problem is to determine the coefficients a, b such that

$$E[(Y - aX - b)^2] \text{ is a minimum}$$

We write the error in a special form, then square and take the expectation.

$$\begin{aligned} \text{Error} &= Y - aX - b = (Y - \mu_Y) - a(X - \mu_X) + \mu_Y - a\mu_X - b = (Y - \mu_Y) - a(X - \mu_X) - \beta \\ \text{Error squared} &= (Y - \mu_Y)^2 + a^2(X - \mu_X)^2 + \beta^2 - 2\beta(Y - \mu_Y) + 2a\beta(X - \mu_X) - 2a(Y - \mu_Y)(X - \mu_X) \\ E[(Y - aX - b)^2] &= \sigma_Y^2 + a^2\sigma_X^2 + \beta^2 - 2a\text{Cov}[X, Y] \end{aligned}$$

Standard procedures for determining a minimum (with respect to a) show that this occurs for

$$a = \frac{\text{Cov}[X, Y]}{\text{Var}[X]} \quad b = \mu_Y - a\mu_X$$

Thus the optimum line, called the *regression line of Y on X* , is

$$u = \frac{\text{Cov}[X, Y]}{\text{Var}[X]}(t - \mu_X) + \mu_Y = \rho \frac{\sigma_Y}{\sigma_X}(t - \mu_X) + \mu_Y = \alpha(t)$$

The second form is commonly used to define the regression line. For certain theoretical purposes, this is the preferred form. But for *calculation*, the first form is usually the more convenient. Only the covariance (which requires both means) and the variance of X are needed. There is no need to determine $\text{Var}[Y]$ or ρ .

Example 12.3.1 The simple air of Example 3 from "Variance"

```
jdemo1
jcalc
Enter JOINT PROBABILITIES (as on the plane)  P
Enter row matrix of VALUES of X  X
Enter row matrix of VALUES of Y  Y
  Use array operations on matrices X, Y, PX, PY, t, u, and P
EX = total(t.*P)
EX =    0.6420
```

```
EY = total(u.*P)
EY =    0.0783
VX = total(t.^2.*P) - EX^2
VX =    3.3016
CV = total(t.*u.*P) - EX*EY
CV =   -0.1633
a = CV/VX
a =   -0.0495
b = EY - a*EX
b =    0.1100           % The regression line is u = -0.0495t + 0.11
```

Example 12.3.2 The pair in Example 6 from "Variance"

Suppose the pair $\{X, Y\}$ has joint density $f_{XY}(t, u) = 3u$ on the triangular region bounded by $u = 0$, $u = 1 + t$, $u = 1 - t$. Determine the regression line of Y on X .

Analytic Solution

By symmetry, $E[X] = E[XY] = 0$, so $\text{Cov}[X, Y] = 0$. The regression curve is

$$u = E[Y] = 3 \int_0^1 u^2 \int_{u-1}^{1-u} dt du = 6 \int_0^1 u^2 (1 - u) du = 1/2$$

Note that the pair is uncorrelated, but by the rectangle test is not independent. With zero values of $E[X]$ and $E[XY]$, the approximation procedure is not very satisfactory unless a very large number of approximation points are employed.

Example 12.3.3 Distribution of Example 5 from "Random Vectors and MATLAB" and Example 12 from "Function of Random Vectors"

The pair $\{X, Y\}$ has joint density $f_{XY}(t, u) = \frac{6}{37}(t + 2u)$ on the region $0 \leq t \leq 2$, $0 \leq u \leq \max\{1, t\}$ (see Figure 12.3.1). Determine the regression line of Y on X . If the value $X(\omega) = 1.7$ is observed, what is the best mean-square linear estimate of $Y(\omega)$?

Figure one contains two lines in the first quadrant of a cartesian graph. The horizontal axis is labeled t, and the vertical axis is labeled u. The title caption reads $f_{xy}(t, u) = (6/37)(t + 2u)$. The first line crosses the vertical axis one quarter of the way up the graph. It has a positive slope, and is labeled $u = 0.3382t + 0.4011$. It continues as a linear plot from one side of the graph to the other. The second line begins horizontally as one segment from the left to point (1, 1). The segment is labeled $u = 1$. After point (1, 1), the line moves upward with a positive, constant slope to point (2, 2). This segment is labeled $u = t$. At (2, 2) there is a vertical line continuing downward to point (2, 0).

Figure 12.3.1. Regression line for Example 12.3.3

Analytic Solution

$$E[X] = \frac{6}{37} \int_0^1 \int_0^1 (t^2 + 2tu) du dt + \frac{6}{37} \int_1^2 \int_0^t (t^2 + 2tu) du dt = 50/37$$

The other quantities involve integrals over the same regions with appropriate integrands, as follows:

Quantity	Integrand	Value
$E[X^2]$	$t^3 + 2t^2u$	779/370
$E[Y]$	$tu + 2u^2$	127/148
$E[XY]$	$t^2u + 2tu^2$	232/185

Then

$$\text{Var}[X] = \frac{779}{370} - \left(\frac{50}{37}\right)^2 = \frac{3823}{13690} \quad \text{and} \quad \text{Cov}[X, Y] = \frac{232}{185} - \frac{50}{37} \cdot \frac{127}{148} = \frac{1293}{13690}$$

and

$$a = \text{Cov}[X, Y] / \text{Var}[X] = \frac{1293}{3823} \approx 0.3382, b = E[Y] - aE[X] = \frac{6133}{15292} \approx 0.4011$$

The regression line is $u = at + b$. If $X(\omega) = 1.7$, the best linear estimate (in the mean square sense) is $\hat{Y}(\omega) = 1.7a + b = 0.9760$ (see Figure 12.3.1 for an approximate plot).

APPROXIMATION

```
tuappr
Enter matrix [a b] of X-range endpoints [0 2]
Enter matrix [c d] of Y-range endpoints [0 2]
Enter number of X approximation points 400
Enter number of Y approximation points 400
Enter expression for joint density (6/37)*(t+2*u).*(u<=max(t,1))
Use array operations on X, Y, PX, PY, t, u, and P
EX = total(t.*P)
EX = 1.3517 % Theoretical = 1.3514
EY = total(u.*P)
EY = 0.8594 % Theoretical = 0.8581
VX = total(t.^2.*P) - EX^2
VX = 0.2790 % Theoretical = 0.2793
CV = total(t.*u.*P) - EX*EY
CV = 0.0947 % Theoretical = 0.0944
a = CV/VX
a = 0.3394 % Theoretical = 0.3382
b = EY - a*EX
b = 0.4006 % Theoretical = 0.4011
y = 1.7*a + b
y = 0.9776 % Theoretical = 0.9760
```

An interpretation of ρ^2

The analysis above shows the minimum mean squared error is given by

$$\begin{aligned} E[(Y - \hat{Y})^2] &= E[(Y - \rho \frac{\sigma_Y}{\sigma_X}(X - \mu_X) - \mu_Y)^2] = \sigma_Y^2 E[(Y^* - \rho X^*)^2] \\ &= \sigma_Y^2 E[(Y^*)^2 - 2\rho X^* Y^* + \rho^2 (X^*)^2] = \sigma_Y^2 (1 - 2\rho^2 + \rho^2) = \sigma_Y^2 (1 - \rho^2) \end{aligned}$$

If $\rho = 0$, then $E[(Y - \hat{Y})^2] = \sigma_Y^2$, the mean squared error in the case of zero linear correlation. Then, ρ^2 is interpreted as the *fraction of uncertainty removed by the linear rule and X*. This interpretation should not be pushed too far, but is a common interpretation, often found in the discussion of observations or experimental results.

More general linear regression

Consider a jointly distributed class. $\{Y, X_1, X_2, \dots, X_n\}$. We wish to determine a function U of the form

$$U = \sum_{i=0}^n a_i X_i, \text{ with } X_0 = 1, \text{ such that } E[(Y - U)^2] \text{ is a minimum}$$

If U satisfies this minimum condition, then $E[(Y - U)V] = 0$, or, equivalently

$$E[YV] = E[UV] \text{ for all } V \text{ of the form } V = \sum_{i=0}^n c_i X_i$$

To see this, set $W = Y - U$ and let $d^2 = E[W^2]$. Now, for any α

$$d^2 \leq E[(W + \alpha V)^2] = d^2 + 2\alpha E[WV] + \alpha^2 E[V^2]$$

If we select the special

$$\alpha = -\frac{E[WW]}{E[V^2]} \text{ then } 0 \leq -\frac{2E[WW]^2}{E[V^2]} + \frac{E[WW]^2}{E[V^2]^2} E[V^2]$$

This implies $E[WW]^2 \leq 0$, which can only be satisfied by $E[WW] = 0$, so that

$$E[UV] = E[UV]$$

On the other hand, if $E[(Y - U)V] = 0$ for all V of the form above, then $E[(Y - U)^2]$ is a minimum. Consider

$$E[(Y - V)^2] = E[(Y - U + U - V)^2] = E[(Y - U)^2] + E[(U - V)^2] + 2E[(Y - U)(U - V)]$$

See $U - V$ is of the same form as V , the last term is zero. The first term is fixed. The second term is nonnegative, with zero value iff $U - V = 0$ a.s. Hence, $E[(Y - V)^2]$ is a minimum when $V = U$.

If we take V to be $1, X_1, X_2, \dots, X_n$, successively, we obtain $n + 1$ linear equations in the $n + 1$ unknowns a_0, a_1, \dots, a_n , as follows.

$$\begin{aligned} E[Y] &= a_0 + a_1 E[X_1] + \dots + a_n E[X_n] \\ E[YX_i] &= a_0 E[X_i] + a_1 E[X_1 X_i] + \dots + a_n E[X_n X_i] \quad \text{for } 1 \leq i \leq n \end{aligned}$$

For each $i = 1, 2, \dots, n$, we take (2) - $E[X_i] \cdot (1)$ and use the calculating expressions for variance and covariance to get

$$\text{Cov}[Y, X_i] = a_1 \text{Cov}[X_1, X_i] + a_2 \text{Cov}[X_2, X_i] + \dots + a_n \text{Cov}[X_n, X_i]$$

These n equations plus equation (1) may be solved algebraically for the a_i .

In the important special case that the X_i are uncorrelated (i.e. $\text{Cov}[X_i, X_j] = 0$ for $i \neq j$), we have

$$a_i = \frac{\text{Cov}[Y, X_i]}{\text{Var}[X_i]} \quad 1 \leq i \leq n$$

and

$$a_0 = E[Y] - a_1 E[X_1] - a_2 E[X_2] - \dots - a_n E[X_n]$$

In particular, this condition holds if the class $\{X_i : 1 \leq i \leq n\}$ is iid as in the case of a simple random sample (see the section on "[Simple Random Samples and Statistics](#)").

Examination shows that for $n = 1$, with $X_1 = X$, $a_0 = b$, and $a_1 = a$, the result agrees with that obtained in the treatment of the regression line, above.

Example 12.3.4 Linear regression with two variables.

Suppose $E[Y] = 3$, $E[X_1] = 2$, $E[X_2] = 3$, $\text{Var}[X_1] = 3$, $\text{Var}[X_2] = 8$, $\text{Cov}[Y, X_1] = 5$, $\text{Cov}[Y, X_2] = 7$, and $\text{Cov}[X_1, X_2] = 1$. Then the three equations are

$$a_0 + 2a_1 + 3a_2 = 3$$

$$0 + 3a_1 + 1a_2 = 5$$

$$0 + 1a_1 + 8a_2 = 7$$

Solution of these simultaneous linear equations with MATLAB gives the results

$a_0 = -1.9565$, $a_1 = 1.4348$, and $a_2 = 0.6957$.

This page titled [12.3: Linear Regression](#) is shared under a [CC BY 3.0](#) license and was authored, remixed, and/or curated by [Paul Pfeiffer](#) via [source content](#) that was edited to the style and standards of the LibreTexts platform.