

EE 464

Spring 2003

Lecture Notes Part 11c

Christopher Wayne Walker, Ph.D.

11.4 Conditional Expected Values

Definition: The *conditional mean of $g(Y)$ given $X \leq x$* is given by

$$E[g(Y)|X \leq x] = \int_{-\infty}^{\infty} g(y)f(y|X \leq x)dy.$$

Definition: The *conditional mean of $g(Y)$ given $X = x$* is given by

$$E[g(Y)|X = x] = \int_{-\infty}^{\infty} g(y)f(y|x)dy.$$

In particular, we have the conditional mean of Y given $X = x$

$$\mu_{Y|X} = E[Y|X = x] = \int_{-\infty}^{\infty} yf(y|x)dy$$

and the conditional variance of Y given $X = x$

$$\sigma_{Y|X}^2 = E[(Y - \mu_{Y|X})^2|X = x] = \int_{-\infty}^{\infty} (y - \mu_{Y|X})^2 f(y|x)dy.$$

Notation: $E[g(Y)|x] = E[g(Y)|X = x]$.

Preceding developments lead to the following theorem.

Theorem:

$$E[g(X, Y)|M] = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} g(x, y)f(x, y|M)dxdy$$

for an event M .

Special Case: Let $M = \{x < X \leq x + \Delta x\}$. Then

$$\begin{aligned} & E[g(X, Y)|x < X \leq x + \Delta x] \\ &= \int_{-\infty}^{\infty} \int_x^{x+\Delta x} g(\alpha, y)f(\alpha, y|x < X \leq x + \Delta x)d\alpha dy. \end{aligned}$$

Recall,

$$f(x, y|x_1 < X \leq x_2) = \frac{f(x, y)}{F_X(x_2) - F_X(x_1)}, \quad x_1 < x \leq x_2.$$

Let $x_1 = x$, $x_2 = x + \Delta x$. Then

$$f(x, y | x_1 < X \leq x + \Delta x) = \frac{f(x, y)}{F_X(x + \Delta x) - F_X(x)}.$$

Therefore,

$$\begin{aligned} & E [g(X, Y) | x < X \leq x + \Delta x] \\ &= \int_{-\infty}^{\infty} \int_x^{x+\Delta x} g(\alpha, y) \frac{f(\alpha, y)}{F_X(\alpha + \Delta x) - F_X(\alpha)} d\alpha dy \\ &= \int_{-\infty}^{\infty} \int_x^{x+\Delta x} g(\alpha, y) f(\alpha, y) \frac{\frac{1}{\Delta x}}{\frac{F_X(\alpha + \Delta x) - F_X(\alpha)}{\Delta x}} d\alpha dy \\ &\rightarrow \int_{-\infty}^{\infty} g(x, y) f(x, y) \Delta x \frac{1}{f_X(x)} dy \quad (\text{as } \Delta x \rightarrow 0). \end{aligned}$$

Thus,

$$E [g(X, Y) | X = x] = \int_{-\infty}^{\infty} g(x, y) \frac{f(x, y)}{f_X(x)} dy$$

which becomes

$$E [g(X, Y) | X = x] = \int_{-\infty}^{\infty} g(x, y) f(y|x) dy.$$

Note that the conditional mean of Y given $X = x$ is itself a function of x :

$$E [Y | X = x] = \int_{-\infty}^{\infty} y f(y|x) dy.$$

Then $E[Y|X]$ is a random variable and

$$\begin{aligned} E [E(Y|X)] &= \int_{-\infty}^{\infty} E(Y|X) f_X(x) dx \\ &= \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} y f(y|x) f_X(x) dy dx. \end{aligned}$$

But,

$$f(y|x) = \frac{f(x, y)}{f_X(x)}$$

so

$$E [E(Y|X)] = \int_{-\infty}^{\infty} y \int_{-\infty}^{\infty} f(x, y) dx dy = \int_{-\infty}^{\infty} y f_Y(y) dy = E(Y).$$

Similarly,

$$E [E (g(X, Y)|X)] = E [g(X, Y)].$$

11.5 Mean Square Estimation

Recall that the value of b that minimizes $E [(X - b)^2]$ is $b = E(X)$ (see class notes section 9.4). So if we wish to estimate the value of a random variable Y using only a constant, c , then the mean square error (MSE)

$$e = E [(Y - c)^2] = \int_{-\infty}^{\infty} (y - c)^2 f_Y(y) dy$$

is minimized if we choose

$$c = E(Y) = \int_{-\infty}^{\infty} y f_Y(y) dy.$$

With $c = E(Y)$, our cost function is $E [(Y - E(Y))^2]$ which is the variance (so we are minimizing the variance in our error).

Nonlinear MS Estimation:

Now consider a possibly nonlinear estimate for Y . Let

$$\begin{aligned} e &= E [(Y - c(X))^2] = \int_{-\infty}^{\infty} (y - c(x))^2 f(x, y) dx dy \\ &= \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} (y - c(x))^2 f(y|x) f_X(x) dx dy \\ &= \int_{-\infty}^{\infty} f_X(x) \int_{-\infty}^{\infty} (y - c(x))^2 f(y|x) dy dx. \end{aligned}$$

Now x is a constant in the integral

$$\int_{-\infty}^{\infty} (y - c(x))^2 f(y|x) dy$$

which implies $c(x)$ is a constant in the integral as well. Since $c(x)$ is a constant we can use our prior result to conclude

$$c(x) = E(Y|X = x) = \int_{-\infty}^{\infty} y f(y|x) dy$$

minimizes this integral for any x . Thus, $E(Y|X = x)$ is the best MSE predictor of Y given $X = x$.

Linear MS Estimation:

Sometimes we are willing to not necessarily have the best minimum mean square estimate (or predictor) but instead a predictor that is easier to calculate.

Theorem: Suppose that $E(X^2)$ and $E(Y^2)$ are finite and X and Y are not constant. Then the best (in the MS sense) zero intercept linear predictor of Y ($\hat{Y} = a_0X$) is obtained by taking

$$a_0 = \frac{E(XY)}{E(X^2)}$$

while the best linear predictor of Y ($\hat{Y} = a_1X + b_1$) is

$$a_1 = \frac{Cov(X, Y)}{Var(X)}, \quad b_1 = E(Y) - a_1E(X).$$

Proof:

$$\begin{aligned} E[(Y - aX)^2] &= E(Y^2) - 2aE(XY) + a^2E(X^2) \\ &= E(X^2) \left[a - \frac{E(XY)}{E(X^2)} \right]^2 + \left[E(Y^2) - \frac{[E(XY)]^2}{E(X^2)} \right]. \end{aligned}$$

Using a we have no control over

$$\left[E(Y^2) - \frac{[E(XY)]^2}{E(X^2)} \right]$$

while

$$E(X^2) \left[a - \frac{E(XY)}{E(X^2)} \right]^2$$

is minimized by taking $a = a_0$. This proves the first part.

Now

$$E[(Y - aX - b)^2] = E[(Y - aX)^2] - 2bE(Y - aX) + b^2$$

$$\begin{aligned}
&= \text{Var}(Y - aX) + [E(Y - aX)]^2 - 2bE(Y - aX) + b^2 \\
&= \text{Var}(Y - aX) + [E(Y)]^2 - 2aE(X)E(Y) + a^2[E(X)]^2 \\
&\quad - 2bE(Y) + 2abE(X) + b^2 \\
&= \text{Var}(Y - aX) + [E(Y) - aE(X) - b]^2.
\end{aligned}$$

Now given any value of a , $[E(Y) - aE(X) - b]^2$ is minimized by taking $b = E(Y) - aE(X) = b_1$. Using this value of b we seek to minimize

$$\begin{aligned}
E[(Y - aX - b)^2] &= E[(Y - aX - (E(Y) - aE(X)))^2] \\
&= E[(Y - E(Y) - a(X - E(X)))^2].
\end{aligned}$$

Let $Y_0 = Y - E(Y)$, $X_0 = X - E(X)$. Then we want to minimize $[(Y_0 - aX_0)^2]$. From first part of theorem we know

$$a = \frac{E(X_0Y_0)}{E(X_0^2)} = \frac{E[(X - E(X))(Y - E(Y))]}{E[(X - E(X))^2]} = \frac{\text{Cov}(X, Y)}{\text{Var}(X)} = a_1.$$

Thus, $b = b_1 = E(Y) - a_1E(X)$ and $\hat{Y} = a_1X + b_1$ is the best linear mean square error predictor or estimator of Y .

Example: Suppose Z_1 is bernoulli with $E(Z_1) = p$ and $\text{Var}(Z_1) = pq$ where $q = 1 - p$. Also, let Z_2 be bernoulli with $E(Z_2) = p$ and $\text{Var}(Z_2) = pq$. Assume Z_1 is independent of Z_2 . Let $X = Z_1$ and let $Y = Z_1Z_2$. Then

- $E(Y|X = x) = E(Z_1Z_2|Z_1 = x) = E(xZ_2) = px$.
- $E[E(Y|X)] = E[pX] = p^2 = E(Y)$.
- $\text{Var}[E(Y|X)] = \text{Var}(pX) = p^2\text{Var}(X) = p^3q$.
- $\text{Var}(Y|X = x) = \text{Var}(Z_1Z_2|Z_1 = x) = \text{Var}(xZ_2) = x^2pq$.
- $E[\text{Var}(Y|X)] = E(X^2pq) = pqE(X^2) = pqE(Z_1^2) = pq[\text{Var}(Z_1) + [E(Z_1)]^2] = pq(pq + p^2) = p^2q(q + p) = p^2q$.
- Best MSE predictor of Y is $E(Y|X) = pX \Rightarrow$ best MSE predictor of Y given $X = x$ is px .

g. Best linear MSE predictor of Y is $\hat{Y} = a_1X + b_1$ where

$$a_1 = \frac{Cov(X, Y)}{Var(X)}, \quad b_1 = E(Y) - a_1E(X).$$

Now

$$\begin{aligned} a_1 &= \frac{E(XY) - E(X)E(Y)}{Var(X)} = \frac{E(Z_1^2Z_2) - E(Z_1)E(Z_1Z_2)}{Var(Z_1)} \\ &= \frac{E(Z_1^2)E(Z_2) - E(Z_1)E(Z_1)E(Z_2)}{Var(Z_1)} \\ &= \frac{(pq + p^2)p - p^3}{pq} = \frac{pq + p^2 - p^2}{q} = p. \end{aligned}$$

So $b_1 = E(Z_1Z_2) - a_1E(Z_1) = E(Z_1)E(Z_2) - a_1E(Z_1) = p^2 - a_1p = p^2 - p^2 = 0$. Thus, $\hat{Y} = pX$.

Therefore the best MSE predictor of Y given X is also the best linear MSE predictor in this case (as expected since the best MSE predictor was itself linear).

Orthogonality Principle

Consider

$$e = E[(Y - (aX + b))^2]$$

where $aX + b$ is a linear estimate of Y given the observed data X . This is minimal where

$$\frac{\partial e}{\partial a} = 0 \quad \text{and} \quad \frac{\partial e}{\partial b} = 0.$$

Thus

$$\frac{\partial e}{\partial b} = E[2(Y - (aX + b))] = 0 \Rightarrow E(Y) = aE(X) + b.$$

Also

$$\frac{\partial e}{\partial a} = E[2(Y - (aX + b))(-X)] = 0 \Rightarrow E[(Y - (aX + b))X] = 0.$$

This implies the estimation error $(Y - (aX + b))$ is orthogonal to the data. This is called the *orthogonality principle*.

Special case: If $b = 0$ we have $e = E[(Y - aX)^2]$ and $E[(Y - aX)X] = 0$ by the orthogonality principle, Thus,

$$E(XY) - aE(X^2) = 0 \Rightarrow a = \frac{E(XY)}{E(X^2)}$$

which is the same as we got for the best zero intercept linear predictor of Y given X .