

Podcast #128: We chat with Kent C Dodds about why he loves React and discuss what life was like in the dark days before Git. [Listen now](#).



Problem with proof of Conditional expectation as best predictor

Asked 6 years, 1 month ago Active 4 months ago Viewed 11k times

▲ I have an issue with the proof of

19

$$E(Y|X) \in \arg \min_{g(X)} E[(Y - g(X))^2]$$

★

which very likely reveal a deeper misunderstanding of expectations and conditional expectations.

16

The proof I know goes as follows (another version of this proof can be found [here](#))

$$\begin{aligned} & \arg \min_{g(X)} E[(Y - g(x))^2] \\ &= \arg \min_{g(X)} E[(Y - E(Y|X) + E(Y|X) - g(X))^2] \\ &= \arg \min_{g(x)} E[(Y - E(Y|X))^2 + 2(Y - E(Y|X))(E(Y|X) - g(X)) + (E(Y|X) - g(X))^2] \\ &= \arg \min_{g(x)} E[2(Y - E(Y|X))(E(Y|X) - g(X)) + (E(Y|X) - g(X))^2] \end{aligned}$$

The proof then typically continues with an argument showing that $2E[(Y - E(Y|X))(E(Y|X) - g(X))] = 0$, and hence

$$\arg \min_{g(x)} E[(Y - g(x))^2] = \arg \min_{g(x)} E[(E(Y|X) - g(X))^2]$$

which can be seen to be minimized when $g(X) = E(Y|X)$.

My puzzles about the proof are the following:

1. Consider

$$E[2(Y - E(Y|X))(E(Y|X) - g(X)) + (E(Y|X) - g(X))^2].$$

It seems to me that, independently of any argument showing that the first term is always equal to zero, one can see that setting $g(X) = E(Y|X)$ minimizes the expression as it implies $(E(Y|X) - g(X)) = 0$ and hence

$$E\left[2(Y - E(Y|X))(E(Y|X) - g(X)) + (E(Y|X) - g(X))^2\right] = E(0 + 0) = 0.$$

But if this is true, then one might repeat the proof replacing $E(Y|X)$ by any other function of X , say $h(X)$, and get to the conclusion that it is $h(X)$ that minimizes the expression. So there must be something I misunderstand (right?).

1. I have some doubts about the meaning of $E[(Y - g(X))^2]$ in the statement of the problem. How should the notation be interpreted? Does it mean

$$E_X[(Y - g(X))^2], E_Y[(Y - g(X))^2] \text{ or } E_{XY}[(Y - g(X))^2]?$$

mathematical-statistics

conditional-probability

proof

conditional-expectation

edited Oct 4 '13 at 23:15

asked Oct 4 '13 at 0:24



Martin Van der Linden

609 1 5 12

4 Answers



(This is an adaptation from Granger & Newbold(1986) "Forecasting Economic Time Series").

11

By construction, your *error cost function* is $[Y - g(X)]^2$. This incorporates a critical assumption (that the error cost function is symmetric around zero) - a different error cost function would not necessarily have the conditional expected value as the $\arg \min$ of its expected value. You cannot minimize your error cost function because it contains unknown quantities. So you decide to minimize its expected value instead. Then your objective function becomes



$$E[Y - g(X)]^2 = \int_{-\infty}^{\infty} [y - g(X)]^2 f_{Y|X}(y|x) dy$$

which I believe answers also your second question. It is intuitive that the expected value will be of Y conditional on X , since we are trying to estimate/forecast Y based on X . Decompose the square to obtain

$$\begin{aligned} E[Y - g(X)]^2 &= \int_{-\infty}^{\infty} y^2 f_{Y|X}(y|x) dy - 2g(X) \int_{-\infty}^{\infty} y f_{Y|X}(y|x) dy \\ &\quad + [g(X)]^2 \int_{-\infty}^{\infty} f_{Y|X}(y|x) dy \end{aligned}$$

The first term does not contain $g(X)$ so it does not affect minimization, and it can be ignored. The integral in the second term equals the conditional expected value of Y given X , and the integral in the last term equals unity. So

$$\arg \min_{g(x)} E[Y - g(X)]^2 = \arg \min_{g(x)} \left\{ -2g(X)E(Y | X) + [g(X)]^2 \right\}$$

The first derivative w.r.t $g(X)$ is $-2E(Y | X) + 2g(X)$ leading to the first order condition for minimization $g(X) = E(Y | X)$ while the second derivative is equal to $2 > 0$ which is sufficient for a minimum.

ADDENDUM: The logic of the "add and subtract" proof approach.

The OP is puzzled by the approach stated in the question, because it seems tautological. It isn't, because while using the tactic of adding and subtracting makes a *specific part* of the objective function zero for an arbitrary choice of the term that is added and subtracted, it does NOT equalize the *value function*, namely the value of the objective function evaluated at the candidate minimizer.

For the choice $g(X) = E(Y | X)$ we have the value function $V(E(Y | X)) = E[(Y - E(Y | X))^2 | X]$ For the arbitrary choice $g(X) = h(X)$ we have the value function $V(h(X)) = E[(Y - h(X))^2 | X]$.

I claim that

$$\begin{aligned} V(E(Y | X)) &\leq V(h(X)) \\ \Rightarrow E(Y^2 | X) - 2E[YE(Y | X) | X] + E[(E(Y | X))^2 | X] \\ &\leq E(Y^2 | X) - 2E[Yh(X) | X] + E[(h(X))^2 | X] \end{aligned}$$

The first term of the LHS and the RHS cancel out. Also note that the outer expectation is conditional on X . By the properties of conditional expectations we end up with

$$\begin{aligned} \dots \Rightarrow -2E(Y | X) \cdot E(Y | X) + [E(Y | X)]^2 &\leq -2E(Y | X)h(X) + [h(X)]^2 \\ \Rightarrow 0 &\leq [E(Y | X)]^2 - 2E(Y | X)h(X) + [h(X)]^2 \\ \Rightarrow 0 &\leq [E(Y | X) - h(X)]^2 \end{aligned}$$

which holds with strict inequality if $h(x) \neq E(Y | X)$. So $E(Y | X)$ is the global and unique minimizer.

But this also says that the "add-and-subtract" approach is not the most illuminating way of proof here.

edited Oct 4 '13 at 10:26

answered Oct 4 '13 at 1:05



Alecos Papadopoulos
45.8k 2 103 208

Thanks for your answer. It helps clarifying my second question. As I tried to convey in the title of the question, my main issue (the first one in the post) was more about the proof mechanism. My main concern is about my understanding of the proof I presented in the question. As I explained, my understanding of the proof leads me to a blatantly problematic statement. So I would like to understand where my mistake is as it might reveal some deeper misunderstandings about concepts of expectation and conditional expectation. Any thoughts about this? – [Martin Van der Linden](#) Oct 4 '13 at 4:00

1 I added some explanation on the "add and subtract" approach to proof. – [Alecos Papadopoulos](#) Oct 4 '13 at 10:27

Took me some time to understand it, but I finally got my basic mistake : true enough $E[-2(Y - h(X))(h(X) - g(X)) + (h(X) - g(X))^2] = 0$ when $g(X) = h(X)$, but by no means does it imply that $h(X)$ minimizes the expression. There is no reason which the bracketed expression could not be lower than zero. Because of the minus sign in front of $(Y - h(X))(h(X) - g(X))$ one could find some $g(X)$ such that $E[-2(Y - h(X))(h(X) - g(X)) + (h(X) - g(X))^2] < 0$. – [Martin Van der Linden](#) Oct 4 '13 at 16:24

1 Hmmm... the minus sign in the expression you refer to is a mistake - it should be a plus sign. You could of course then rearrange the terms to obtain again a minus sign... does this hurt the intuition you gained? – [Alecos Papadopoulos](#) Oct 4 '13 at 21:03

Thanks for keeping up with the question. I edited the initial post to correct for this mistake. Fortunately, I think it does not hurt the gained intuition. Actually it helps me understand yet another mistake : I was assuming that the minus sign was important to guarantee that 0 was not necessarily the minimum of $E[-2(Y - h(X))(h(X) - g(X)) + (h(X) - g(X))^2]$. But I realize this is not just about the sign before the 2. (Hopefully) What I really needed to understand is that, in general (i.e. for arbitrary $h(X)$) $E[2(Y - h(X))(h(X) - g(X))]$ needs not be minimized when $g(X) = h(X)$ (right?). – [Martin Van der Linden](#) Oct 4 '13 at 23:36

|

Note that to prove the answer, you really only need to show that

5

$$E[-2(Y - E(Y|X))(E(Y|X) - g(X))] = 0$$

As for which expectation to take, you take it conditionally, otherwise the term

$$\arg \min_{g(X)} E[(Y - g(X))^2]$$

Doesn't make sense, as $g(X)$ is a random variable if E is E_{XY} and not $E_{Y|X}$. Show you should really write $E[(Y - g(X))^2|X]$ or $E_{Y|X}[(Y - g(X))^2]$ to make this clear. Now given this clarification, the term $(E(Y|X) - g(X))$ is a constant, and can be pulled outside the expectation, and you have:

$$\begin{aligned} & -2(E(Y|X) - g(X))E[(Y - E(Y|X))|X] = \\ & -2(E(Y|X) - g(X))[E(Y|X) - E[E(Y|X)|X]] = \\ & -2(E(Y|X) - g(X))[E(Y|X) - E(Y|X)] = 0 \end{aligned}$$

Hence you can write the objective function as:

$$E_{Y|X} \left[(Y - g(X))^2 \right] = E_{Y|X} \left[(Y - E_{Y|X}(Y|X))^2 \right] + (E_{Y|X}(Y|X) - g(X))^2$$

The minimiser is obvious from here. Note that if you were to average over X as well, then a very similar argument can be used to show:

$$\begin{aligned} E_X \left[(E(Y|X) - g(X))^2 \right] &= E_X \left[(E_{Y|X}(Y|X) - E_X[E_{Y|X}(Y|X)])^2 \right] \\ &\quad + \left(E_X[E_{Y|X}(Y|X)] - E_X[g(X)] \right)^2 \end{aligned}$$

This shows that if you set $g(X) = E_{Y|X}(Y|X)$ for each X , then you also have a minimiser over this function as well. So in some sense it doesn't really matter whether E is E_{YX} or $E_{Y|X}$.

answered Oct 4 '13 at 4:57



probabilityislogic

20.6k 3 66 89



There's a mathematical point of view that is very simple. What you have is a projection problem in a Hilbert space, much like projecting a vector in \mathbb{R}^n onto a subspace.

3



Let $(\Omega, \mathcal{F}, \mu)$ denote the underlying probability space. For the problem to make sense, consider the random variables with finite second moments, that is, the Hilbert space $L^2(\Omega, \mathcal{F}, \mu)$. The problem now is this: given $X, Y \in L^2(\Omega, \mathcal{F}, \mu)$, find the projection of Y onto the subspace $L^2(\Omega, \mathcal{F}_X, \mu)$, where \mathcal{F}_X is the σ -subalgebra of \mathcal{F} generated by X . (Just as in the finite dimensional case, minimizing L^2 -distance to a subspace means finding the projection). The desired projection is $E(X|Y)$, by construction. (This actually characterizes $E(X|Y)$, if one inspects the proof of existence).

answered Nov 3 '13 at 3:17



Michael

658 3 10



This is a beautiful response. – **jll** Mar 18 '15 at 20:36



Regarding your last question, the expectation can be either w.r.t. $p(x, y)$ (the unconditional error) or w.r.t. $p(y | x)$ (the conditional error at each value $X = x$). Happily, minimizing the conditional error at each value $X = x$ also minimizes the unconditional error, so this is not a crucial distinction.

0



answered Jun 25 at 2:36



Ulisses Braga-Neto

101 1