

6.5 Conditional Expectation

There are 2 closely linked notions of conditional expectation:

- Conditional expectation $E(Y|A)$ given an event $\{A\}$: let Y be an r.v., and A be an event. If we learn that A occurred, our updated expectation for Y , $E(Y|A)$, is computed analogously to $E(Y)$, except using conditional probabilities given A .
- Conditional expectation $E(Y|X)$ given an r.v.: intuitively, $E(Y|X)$ is the r.v. that best predicts Y using only information from X .

Definition 6.5.1 (Conditional expectation given an event)

Let A be event with positive probability. The conditional expectation of Y given A is:

- If Y is discrete r.v.: $E(Y|A) = \sum_y y \cdot P(Y=y|A)$
- If Y is continuous r.v.: $E(Y|A) = \int_{-\infty}^{\infty} y \cdot f(y|A) dy$,
where $f(y|A)$ is the derivative of CDF $F(y|A) = P(Y \leq y|A)$,
and can be computed using Bayes' rule:

$$f(y|A) = \frac{P(A|Y=y) f(y)}{P(A)}$$

Warning 6.5.2 (Mistaken conditional and unconditional expectation)

Confusing conditional and unconditional expectation is dangerous. You should always be aware of what you are conditioning on.

Example 6.5.3 (Life expectancy)

Fred is 30 years old, his country life expectancy is 80. Should he conclude that, on average he has 50 years left to live?

No, because he missed 1 crucial fact: he has already lived to 30.

Let T be Fred's lifespan:

$$E(T) < E(T \mid T \geq 30)$$

To prove this, consider:

$$\begin{aligned} E(T \mid T \geq 30) &= \sum_{t=30}^{\infty} t \cdot P(T=t \mid T \geq 30) \\ &= \sum_{t=30}^{\infty} t \cdot \frac{P(T=t)}{P(T \geq 30)} \\ &= \frac{1}{P(T \geq 30)} \sum_{t=30}^{\infty} t \cdot P(T=t) \\ &\geq 1 \end{aligned}$$

subset of $E(T) = \sum_{t=0}^{\infty} t \cdot P(T=t)$

Theorem 6.5.4 (Law of total expectation)

Let A_1, \dots, A_n be a partition of a sample space, with $P(A_i) > 0$ for all i , and let Y be a random variable on this sample space.

Then:

$$E(Y) = \sum_{i=1}^n E(Y \mid A_i) P(A_i)$$

Connection with LOTP through fundamental bridge:

Let $Y = I_B$ for an event B , then:

$$P(B) = E(I_B) = \sum_{i=1}^n E(I_B \mid A_i) P(A_i) = \sum_{i=1}^n P(B \mid A_i) P(A_i)$$

Example 6.5.5 (Two-envelope paradox)

There are 2 sealed envelopes, one of which has twice the amount of money to the other. Which one should you choose?

Let X and Y represent the amount in the left and right envelopes.

. There are 2 theories:

- By symmetry that $E(X) = E(Y)$, there is no reason that you should choose one over the other
- Suppose that the left has \$100, then the right envelope has either \$50 or \$200, which average to \$125. This contradict the first theory, and since we can apply the same reasoning to the right envelope, leading to switching back and forth between 2 envelopes.

Formalizing this, we have $Y = 2X$ or $Y = \frac{X}{2}$ with equal probabilities.

By LOTE:

$$E(Y) = E(Y | Y=2X) \cdot \frac{1}{2} + E(Y | Y=\frac{X}{2}) \cdot \frac{1}{2}$$

One might think that:

$$E(Y) = E(2X) \cdot \frac{1}{2} + E\left(\frac{X}{2}\right) \cdot \frac{1}{2} = \frac{5}{4} E(X)$$

which is 25% from switching (by theory 2).

However, there is no justification for $E(Y | Y=2X) = E(2X)$ or $E(Y | Y=\frac{X}{2}) = E\left(\frac{X}{2}\right)$. To see why, let I indicate $Y=2X$, then:

$$E(Y | 2X) = E(2X | I=1)$$

(we can only drop this term if X and I are independent, but they are not (definition of I above))

Example 6.5.7 (Mystery prize)

You are bidding on a mystery box containing a mystery prize. The prize is unknown, except that its worth between 0-1 million dollars. So the true value V of the prize is Uniform on $[0, 1]$ (in million of dollars)

You can bid any amount b (in millions of dollars). Specifically, the bid is

| |
|--|
| } rejected if $b < \frac{2}{3}V$ |
| } accepted if $b \geq \frac{2}{3}V$, your payoff is $V-b$ |

How much should you bid?

Let W be the payoff amount.

To find the expected payoff amount, condition on whether the bid is accepted. The payoff is

| |
|--------------------------------|
| } $V-b$ if the bid is accepted |
| } 0 if the bid is rejected |

So by LOTE:

$$\begin{aligned} E(W) &= E(W | b \geq \frac{2}{3}V) \cdot P(b \geq \frac{2}{3}V) + E(W | b < \frac{2}{3}V) \cdot P(b < \frac{2}{3}V) \\ &= E(V-b | b \geq \frac{2}{3}V) \cdot P(b \geq \frac{2}{3}V) + E(0 | b < \frac{2}{3}V) \cdot P(b < \frac{2}{3}V) \\ &= E(V-b | b \geq \frac{2}{3}V) \cdot P(b \geq \frac{2}{3}V) \quad \text{Rejected case} \\ &= (E(V | V \leq \frac{3}{2}b) - b) \cdot P(V \leq \frac{3}{2}b) \quad = 0 \end{aligned}$$

For $b \geq \frac{2}{3}V$, $P(V \leq \frac{3}{2}b) = P(V \leq 1) = 1 \quad \leftarrow V \sim \text{Uniform}$

So the RHS equals $\frac{1}{2} - b < 0 \Rightarrow$ you lose money

For $b < \frac{2}{3}V$, $P(V \leq \frac{3}{2}b) = \frac{3}{2}b$, so the conditional distribution of V is Uniform on $[0, \frac{3}{2}b]$.

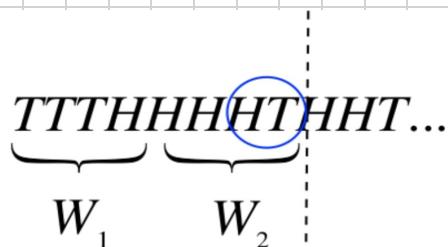
$$\begin{aligned} \text{So the RHS equals } & (E(V | V \leq \frac{3}{2}b) - b) \cdot P(V \leq \frac{3}{2}b) \\ &= \left(\frac{3}{4}b - b\right) \left(\frac{3}{2}b\right) \\ &= -\frac{3}{8}b^2 \Rightarrow \text{You also lose money} \end{aligned}$$

In conclusion, you should not play this game.

Example 6.5.8 (Time until HT vs HT)

A fair coin is tossed repeatedly. What is the expected number of tosses until the pattern HT appears for the first time? What is the expected number of tosses until the pattern HH appears for the first time?

- Let W_{HT} be the number of tosses until HT appears.



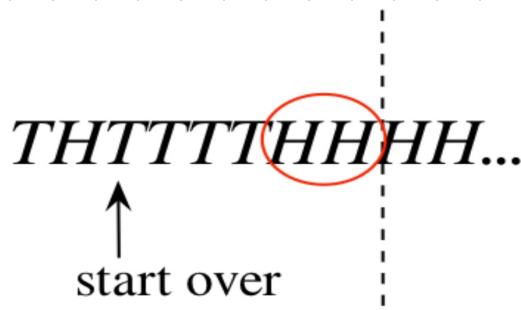
From the image we can see that $W_{HT} = W_1 + W_2$

so by linearity, $E(W_{HT}) = E(W_1) + E(W_2)$

And since $W_1, W_2 \stackrel{\text{i.i.d}}{\sim} \text{FS}\left(\frac{1}{2}\right)$ (W_1 is tosses until first Head,
 W_2 is tosses until first Tail)

$$\Rightarrow E(W_1) = E(W_2) = \frac{1}{p} = 2 \Rightarrow E(W_{HT}) = 4$$

- Let W_{HH} be number of tosses until HH appears.



From the image if we toss
 } . Head, then we progress to the next stage
 } and wait for another head
 } . Tail, then the problem reset

The "reset" nature of the problem makes it perfect for solving by conditioning on the first step, so by LOTE:

$$E(W_{HH}) = E(W_{HH} \mid \text{first toss H}) \cdot \underbrace{\frac{1}{2}}_{\text{first term}} + E(W_{HH} \mid \text{first toss T}) \cdot \underbrace{\frac{1}{2}}_{\text{second term}}$$

- For second term, $E(W_{HH} | \text{first toss } T) = 1 + E(W_{HH})$ by memorylessness.
- For first term, $E(W_{HH} | \text{first toss } H)$ is computed by further conditioning on the second toss, so:

$$\begin{aligned}
 E(W_{HH} | \text{first toss } H) &= E(W_{HH} | \text{first toss } H, \text{second toss } H) \cdot \frac{1}{2} \\
 &\quad + E(W_{HH} | \text{first toss } H, \text{second toss } T) \cdot \frac{1}{2} \\
 &= \underbrace{2 \cdot \frac{1}{2}}_{\substack{\text{get pattern} \\ \text{HH in 2 tosses}}} + \underbrace{(2 + E(W_{HH})) \cdot \frac{1}{2}}_{\substack{\text{wasted 2 tosses} \\ \text{and start over}}}
 \end{aligned}$$

- Replacing first and second terms in the original equation, gives:

$$E(W_{HTH}) = \left(2 \cdot \frac{1}{2} + (2 + E(W_{HH})) \cdot \frac{1}{2}\right) \cdot \frac{1}{2} + (1 + E(W_{HH})) \cdot \frac{1}{2}$$

Solve for $E(W_{HH})$, we get $E(W_{HH}) = 6$

Why $E(W_{HH}) = 6 > E(W_{HT}) = 4$ when probability of event HH and HT are both $\frac{1}{4}$?

- When waiting for HT, if H is tossed, we achieve partial progress where if we get another H, we are at the same position and if we get T, we are done.
- When waiting for HH, if H is tossed and follows by T, we lose all our progress, this suggest average waiting time for HH is longer.
- Symmetry implies that average waiting time for HH and TT are the same, and HT and TH is the same, not HH and HT.

Definition 6.6.1 (Conditional Expectation given a random variable)

- Let $g(x) = E(Y | X=x)$. Then the conditional expectation of Y given X , denoted $E(Y | X)$, is defined to be random variable $g(X)$.
- In other words, if after doing experiment X crystallizes into x , then $E(Y | X)$ crystallizes into $g(x)$.
- If Y is discrete:

$$g(x) = E(Y | X=x) = \sum_y y \cdot P(Y=y | X=x)$$

- If Y is continuous:

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

Warning 6.6.2 (Notion confusion)

The notion can sometimes cause confusion. It doesn't say:

$$g(x) = E(Y | X=x), \text{ so } g(X) = E(Y | X=X) = E(Y)$$

Rather, we should first compute function $g(x)$, then plug X for x .

For example, if $g(x) = x^2$, then $g(X) = X^2$

Warning 6.6.3 (Difference between $E(Y|X)$ and $E(Y|A)$)

- By definition, $E(Y|X)$ is a function of X , so it's a random variable.
Thus, make sense to compute values like the mean $E(E(Y|X))$,
and variance $\text{Var}(E(Y|X))$.
- Worth keeping in mind that:
 - $E(Y|A)$ are numbers
 - $E(Y|X)$ are random variables

Example 6.6.4 (Mean and variance of $E(Y|X)$)

Suppose we have a stick of length 1, we break it at point X chosen uniformly at random. Given that $X=x$, we then choose another break point Y uniformly on interval $[0, x]$.

Find $E(Y|X)$, and its mean, variance.

. We have $X \sim \text{Unit}(0, 1)$ and $Y|X=x \sim \text{Unit}(0, x)$, so:

$$E(Y|X=x) = \frac{x}{2}$$

Plugging x for x , we have:

$$E(Y|X) = \frac{X}{2}$$

. The expected value of $E(Y|X)$ is:

$$E(E(Y|X)) = E\left(\frac{X}{2}\right) = \frac{1}{4}$$

. The variance of $E(Y|X)$ is:

$$\text{Var}(E(Y|X)) = \text{Var}\left(\frac{X}{2}\right) = \frac{1}{48}$$

6.7 Adam's law and other properties of conditional expectation

Some useful properties of conditional expectation

Theorem 6.7.1 (Dropping what's independent)

If X and Y are independent, then $E(Y|X) = E(Y)$

Proof:

- Independence implies $E(Y|X=x) = E(Y)$ for all x , so $E(Y|X) = E(Y)$
- Intuitively, if X provides no information about Y , the our best guess for Y , even if we know X , is still $E(Y)$.

The converse is false:

Counterexample is given in Example 6.7.4

Example 6.7.4 ($E(Y|X) = E(Y)$ doesn't mean X and Y are independent)

- Let X_1, \dots, X_n be i.i.d., and $S_n = X_1 + \dots + X_n$. Find $E(X_i|S_n)$
- By linearity: $E(X_1|S_n) + \dots + E(X_n|S_n) = E(S_n|S_n) = S_n$
- And by symmetry: $E(X_1|S_n) = \dots = E(X_n|S_n)$
- Therefore: $E(X_i|S_n) = \frac{S_n}{n} = \bar{X}_n$
which is the sample mean of all X_j 's

Theorem 6.7.2 (Taking out what's known)

For any function h ,

$$E(h(x). Y | X) = h(x) \cdot E(Y | X)$$

Proof:

Intuitively, when we take expectations condition on X , we are treating as if X crystallize into x . Then any function of X , like $h(X)$, also crystallize into $h(x)$.

Theorem 6.7.3 (Linearity)

$$E(Y_1 + Y_2 | X) = E(Y_1 | X) + E(Y_2 | X)$$

Proof:

This is arrived from the fact that: $E(Y_1 + Y_2) = E(Y_1) + E(Y_2)$

Theorem 6.7.4 (Adam's Law)

For any r.v.s X and Y ,

$$E(E(Y | X)) = E(Y)$$

Proof:

This is proof for X and Y as discrete

Let $E(Y | X) = g(X)$, we have

$$E(X) = \sum_x x \cdot P(X=x)$$

< definition of expectation >

$$\Leftrightarrow E(g(X)) = \sum_x g(x) \cdot P(X=x)$$

< LOTUS >

$$= \sum_x \left(\sum_y y \cdot P(Y=y | X=x) \right) \cdot P(X=x) \quad < \text{definition of } g(x) >$$

$$= \sum_x \sum_y y \cdot P(X=x) \cdot P(Y=y | X=x)$$

$$= \sum_x \sum_y y \cdot P(X=x, Y=y) \quad < \text{Bayes' rule} >$$

$$= \sum_y y \cdot \sum_x P(X=x, Y=y) \quad < \text{Swapping sum} >$$

$$= \sum_y y \cdot P(Y=y) \quad < \text{sum over all } x >$$

$$= E(Y) \quad < \text{definition of expectation of } Y >$$

o Apply the same logic for continuous case.

o Adam's Law connects conditional expectation to unconditional expectation

Example 6.7.6 (Linear regression)

linear regression predict Y based off X , it assumes the conditional expectation of Y is linear in X :

$$E(Y|X) = a + bX$$

a) Show that this is equivalent to: $Y = a + bX + \epsilon$,

where ϵ is a r.v (error) with $E(\epsilon|X) = 0$.

we have: $Y = a + bX + \epsilon$, where $E(\epsilon|X) = 0$

$$\begin{aligned} \Leftrightarrow E(Y|X) &= E(a|X) + E(bX|X) + E(\epsilon|X) \quad \langle \text{linearity} \rangle \\ &= E(a) + bX E(1|X) + 0 \quad \langle \text{properties of conditional} \rangle \\ &= a + bX \end{aligned}$$

b) Solve for constants a and b in terms of $E(X)$, $E(Y)$, $\text{Cov}(X, Y)$ and $\text{Var}(X)$

From a: $E(Y|X) = a + bX + \epsilon$

$$\Leftrightarrow E(E(Y|X)) = a + bE(X) + E(\epsilon) \quad \langle \text{Expectation both sides} \rangle$$

$$\Leftrightarrow E(Y) = a + bE(X) + 0 \quad \left\{ \text{LHS: Adam's Law, RHS: } E(\epsilon) = E(E(\epsilon|X)) = 0 \right\}$$

$$\Leftrightarrow a = E(Y) - bE(X)$$

Taking covariance with X of both sides in $Y = a + bX + \epsilon$:

$$\begin{aligned} \text{Cov}(X, Y) &= \text{Cov}(X, a) + b\text{Cov}(X, X) + \text{Cov}(X, \epsilon) \\ &= b \text{Var}(X) \end{aligned}$$

Thus:
$$\left\{ \begin{array}{l} b = \frac{\text{Cov}(X, Y)}{\text{Var}(X)} \\ a = E(Y) - \frac{\text{Cov}(X, Y)}{\text{Var}(X)} E(X) \end{array} \right.$$

6.8 Eve's law and conditional variance

Definition 6.8.1 (Conditional Variance)

The conditional variance of Y given X is:

$$\text{Var}(Y|X) = E((Y - E(Y|X))^2 | X)$$

This is equivalent to:

$$\text{Var}(Y|X) = E(Y^2|X) - (E(Y|X))^2$$

→ Almost the same as unconditional variance:

We basically just replace all instances of $E(\cdot)$ in the definition with $E(\cdot|X)$

→ Connection with conditional expectation

Since conditional variance is defined by conditional expectations, we can use results about conditional expectations to calculate conditional variance.

Warning 6.8.2 (conditional variance is a random variable)

like $E(Y|X)$, $\text{Var}(Y|X)$ is a random variable, and is a function of X .

Theorem 6.8.4 (Eve's Law) or Law of Total Variance

For any r.v.s X and Y ,

$$\text{Var}(Y) = E(\text{Var}(Y|X)) + \text{Var}(E(Y|X))$$

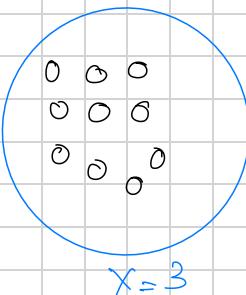
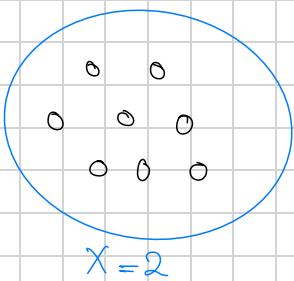
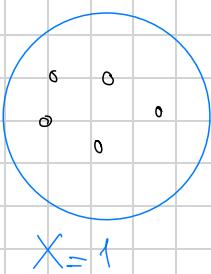
→ Proof: Let $g(X) = E(Y|X)$. By Adam's Law, $E(g(X)) = E(Y)$.

$$\begin{aligned} E(\text{Var}(Y|X)) &= E(E(Y^2|X) - g(X)^2) = E(Y^2) - E(g(X)^2) \\ \text{Var}(E(Y|X)) &= E(g(X)^2) - (Eg(X))^2 = E(g(X)^2) - (EY)^2 \end{aligned}$$

Adding these we have $\text{Var}(Y)$

• Visualizing Eve's Law

We split sample space by r.v X , i.e:



Let Y be r.v represent number of pebbles

Eve's Law says, in general:

$$\text{total variance} = \text{within-group variance} + \text{between-group variance}$$

$$\begin{aligned} E(Y) & \quad (\text{average pebbles within} \\ & \quad \text{each group}) & & \quad (\text{variation of group} \\ & \quad \text{means}) \\ E(\text{Var}(Y|X)) & & & \quad \text{Var}(E(Y|X)) \end{aligned}$$

Warning 6.8.6 (Confusion with LOTE)

It's wrong to say that:

$$\text{Var}(Y) = \text{Var}(Y|A) \cdot P(A) + \text{Var}(Y|A^c) \cdot P(A^c)$$

Rather, we should use Eve's Law, let I be indicator event A :

$$\text{Var}(Y) = E(\text{Var}(Y|I)) + \text{Var}(E(Y|I))$$

Example 6.8.7 (Random sum)

A store has N customers in a day, where N is r.v with finite mean and variance.

Let X_j be the amount customer j^{th} spend in that day, assume X_j has mean μ and variance σ^2

N and all of X_j are independent of each other.

Find the mean and variance of random sum $X = \sum_{j=1}^N X_j$, which is the store revenue in a day, in terms of μ , σ^2 , $E(N)$ and $\text{Var}(N)$

- We might try to find mean by linearity: $E(X) = N\mu$

But: $\left. \begin{array}{l} E(X) \text{ is function of } N \text{ and } X_j, X \text{ is not a sum (number), but a} \\ \text{random sum (function of } N \text{ and } X_j) \end{array} \right\} N\mu \text{ is a random variable (function of } N)$

And linearity only applies to fixed numbers (constants)

- So let's X to function of just N , by condition on N :

$$\begin{aligned} E(X|N) &= E\left(\sum_{j=1}^N X_j | N\right) && \langle \text{Linearity of conditional expectation} \rangle \\ &= \sum_{j=1}^N E(X_j | N) && \langle \text{Taking out what's known} \rangle \\ &= \sum_{j=1}^N E(X_j) && \langle X_j \text{ and } N \text{ are independent} \rangle \\ &= N \cdot \mu \end{aligned}$$

" $E(X|N) = N\mu$ " is valid since both sides of the equation are random variables (functions of N)

- Now we can find $E(X)$ by Adam's Law:

$$E(X) = E(E(X|N)) = E(N\mu) = \mu E(N)$$

This means:

Average Revenue = average spend per customer \times average number of customers

- To find $\text{Var}(X)$, we condition on N to turn $\text{Var}(X)$ from function of N and X_j (since X is function of N and X_j) to function of just N .

$$\begin{aligned} \text{Var}(X|N) &= \text{Var}\left(\sum_{j=1}^N X_j | N\right) \\ &= \sum_{j=1}^N \text{Var}(X_j | N) && \langle \text{taking out what's known} \rangle \\ &= \sum_{j=1}^N \text{Var}(X_j) = N\sigma^2 && \langle X_j \text{ and } N \text{ independent} \rangle \end{aligned}$$

Then use Eve's law to get $\text{Var}(X)$:

$$\begin{aligned}\text{Var}(X) &= E(\text{Var}(X|N)) + \text{Var}(E(X|N)) \\ &= E(N\sigma^2) + \text{Var}(N\mu) \\ &= \sigma^2 E(N) + \mu \text{Var}(N)\end{aligned}$$