

# MENTOR GUIDE TO LLSE 1

## COMPUTER SCIENCE MENTORS 70

November 14 to November 18, 2016

### 1 Covariance

#### 1.1 Formulas

Provided that  $X$  and  $Y$  are random variables, we define **covariance** as follows:

$$\text{Cov}(X, Y) = E((X - E(X)) \cdot (Y - E(Y)))$$

**Claim:**  $\text{Cov}(X, Y) = E(XY) - E(X)E(Y)$

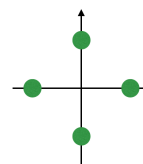
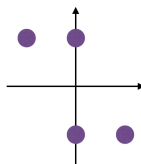
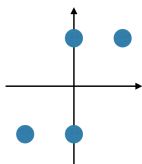
$$\begin{aligned} E((X - E(X)) \cdot (Y - E(Y))) &= E(XY - E(X)Y - XE(Y) + E(X)E(Y)) \\ &= E(XY) - E(X)E(Y) - E(X)E(Y) + E(X)E(Y) \end{aligned}$$

We can cancel out the last two terms:

$$= E(XY) - E(X)E(Y)$$

#### 1.2 Intuition

Look at the following graphs. What is the covariance of each? If  $X$  increases, what happens to  $Y$ ? Are these events independent? Does uncorrelated imply independence? Does independence imply uncorrelated?



### 1.3 Properties of Covariance

---

1.  $\text{Var}(X) = \text{Cov}(X, X)$
2.  $X, Y$  independent  $\rightarrow \text{Cov}(X, Y) = 0$ . Independent random variables are uncorrelated. The converse is not necessarily true (look at the last example on the previous page).
3.  $\text{Cov}(a + X, b + Y) = \text{Cov}(X, Y)$
4.  $\text{Cov}(aX + bY, cU + dV) = ac \cdot \text{Cov}(X, U) + ad \cdot \text{Cov}(X, V) + bc \cdot \text{Cov}(Y, U) + bd \cdot \text{Cov}(Y, V)$

#### Proofs

1. Plug into original equation.
2. Plug into original equation (see what you can cancel out by applying the definition of covariance and independence of the variables).
3. Plug into original equation.
4. This one is a little tricky. Since we know that we can add and subtract **constants**, we can zero mean the variables. So we can assume that  $E(aX + bY) = E(cU + dV) = 0$ . Now let's apply the definition of covariance.

$$\text{Cov}(aX + bY, cU + dV)$$

Since  $E(aX + bY) = E(cU + dV) = 0$ , we get

$$= E((aX + bY) \cdot (cU + dV)) - 0 \cdot 0$$

By linearity of expectation:

$$= ac \cdot E(XU) + ad \cdot E(XV) + bc \cdot E(YU) + bd \cdot E(YV)$$

Now note that  $\text{Cov}(X, U) = E(XU) - E(X)E(U) = E(XU) - 0 \cdot 0 = E(XU)$

Apply the above result to each term to get:

$$= ac \cdot \text{Cov}(XU) + ad \cdot \text{Cov}(XV) + bc \cdot \text{Cov}(YU) + bd \cdot \text{Cov}(YV)$$

**2 LLSE****2.1 Constant**

**Motivation:** Say we have a random variable  $Y$ , which measures height. What would be our best guess for  $Y$ ? In other words, if a random person called us, what would be our best guess for their height?

**2.2 Minimizing Error**

We need to first somehow quantify "the best guess". When we choose a number for height, how do we know that the number we chose is a good guess? It makes sense to pick such a number, where the error is minimized. We will look at the **mean square error**,

$$E((Y - a)^2)$$

and show that this equation achieves a minimum value only when  $a = E(Y)$ .

Define

$$\hat{Y} = Y - E(Y)$$

$\hat{Y}$  is the error we make by using  $E(Y)$  to estimate  $Y$ .

Note that,

$$E(\hat{Y}) = E(Y - E(Y)) = E(Y) - E(E(Y)) = E(Y) - E(Y) = 0$$

This implies that  $E(\hat{Y} \cdot c) = 0 \forall c$

Now we will use these facts to show that if  $a$  is any constant other than  $E(Y)$ , the mean square error will be larger than if we had used  $E(Y)$ .

$E((Y - a)^2) = E((Y - E(Y) + E(Y) - a)^2)$	trick
$= E((\hat{Y} + c)^2)$	where $c = E(Y) - a$
$= E(\hat{Y}^2 + 2\hat{Y}c + c^2)$	expand the square
$= E(\hat{Y}^2) + 2E(\hat{Y}c) + c^2$	linearity of expectation
$= E(\hat{Y}^2) + 0 + c^2$	using fact from above
$\geq E(\hat{Y}^2)$	$c^2 \geq 0$

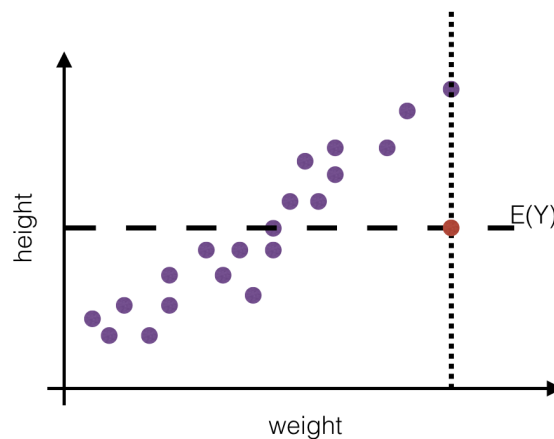
Hence,

$$E((Y - a)^2) \geq E((Y - E(Y))^2)$$

## 2.3 Linear

We now have a good estimate for  $Y$  if we have no additional information. However, what if we know some samples of  $Y$  and some other random variable  $X$ . For example, say we know the weight of the caller. We should be better able to guess their height, as there is some sort of correlation between weight and height. However, the model we are using so far does not apply any additional information to the guess. Let us construct a new way to predict  $Y$ , based on some observed  $X$ .

Continuing the example we have been using so far with height and weight, assume we are given the the points shown below.



We can see that  $E(Y)$  is not a good guess given additional information. We wish to construct a linear equation which given a value of  $X$ , will predict the value of  $Y$ . So if someone calls you and provides their weight, you can more accurately predict what their height is.

So far, our example has been non Bayesian, meaning we did not know the distributions of  $X$  and  $Y$ . We will develop the linear least square estimate using  $X$ ,  $Y$  with known distributions. Then we will see how we can apply these results to non Bayesian stuff.

**Theorem** Assume  $X, Y$  are random variables with known distribution. Then,

$$L[Y|X] = E(Y) + \frac{\text{Cov}(X, Y)}{\text{Var}(X)} \cdot (X - E(X)) = \hat{Y}$$

We want to show that the mean square between  $Y$  and  $\hat{Y}$  is smaller than between  $Y$  and any other linear function of  $X$ . We will use a technique similar to how we showed the constant case.

**2.4 First Proof of  $L[Y|X]$** 

---

Recall,

$$\hat{Y} = E(Y) + \frac{\text{Cov}(X, Y)}{\text{Var}(X)} \cdot (X - E(X))$$

Plug the above equation into the mean value of the error:

$$\begin{aligned} Y - \hat{Y} &= Y - E(Y) - \frac{\text{Cov}(X, Y)}{\text{Var}(X)} \cdot (X - E(X)) \\ E(Y - \hat{Y}) &= E\left(Y - E(Y) - \frac{\text{Cov}(X, Y)}{\text{Var}(X)} \cdot (X - E(X))\right) \\ &= E(Y - E(Y)) - \frac{\text{Cov}(X, Y)}{\text{Var}(X)} E(X - E(X)) \\ &= 0 - 0 = 0 \end{aligned}$$

So we know that the mean value of the error is 0.

**Claim:**  $E((Y - \hat{Y})X) = 0$

**Proof of claim:**

Note,

$$E((Y - \hat{Y})X) = E((Y - \hat{Y}) \cdot (X - E(X)))$$

To see that the two sides are equivalent we do the following: expand the LHS and apply linearity of expectation we get,

$$\begin{aligned} E((Y - \hat{Y}) \cdot (X - E(X))) &= E(YX - \hat{Y}X - YE(X) + \hat{Y}E(X)) \\ &= E(YX - \hat{Y}X) - E(X)E(Y - \hat{Y}) \quad \text{mean value of error is 0} \\ &= E(YX - \hat{Y}X) = E((Y - \hat{Y})X) \end{aligned}$$

Now simplify the new equation,  $E((Y - \hat{Y}) \cdot (X - E(X)))$ :

$$\begin{aligned} E((Y - \hat{Y}) \cdot (X - E(X))) &= E\left(\left(Y - \left(E(Y) + \frac{\text{Cov}(X, Y)}{\text{Var}(X)} \cdot (X - E(X))\right)\right) \cdot (X - E(X))\right) \\ &= E((Y - E(Y))(X - E(X))) - \frac{\text{Cov}(X, Y)}{\text{Var}(X)} E((X - E(X))(X - E(X))) \\ &= \text{Cov}(X, Y) - \frac{\text{Cov}(X, Y)}{\text{Var}(X)} \text{Var}(X) \\ &= 0 \end{aligned}$$

We conclude that  $E((Y - \hat{Y})X) = 0$ .

We need one more fact in order to prove our original claim:

**Claim:** The error is orthogonal to all linear functions. (We have shown that it is orthogonal to  $x$ )

**Proof of claim:**

$$\begin{aligned} E((Y - \hat{Y})(c + dX)) &= E(Yc - \hat{Y}c + YdX - \hat{Y}dX) \\ &= cE(Y - \hat{Y}) + dE((Y - \hat{Y})X) \\ &= 0 - 0 = 0 \end{aligned} \quad \text{using above results}$$

Now let's get back to our original claim: We want to show that the mean square between  $\hat{Y}$  and  $Y$  is smaller than between  $Y$  and any other linear function. Recall that to show  $E(Y)$  is the best constant estimate for  $Y$  we proved the following equation:

$$E((Y - a)^2) \geq E((Y - E(Y))^2)$$

To show that  $\hat{Y}$  is the best linear estimator for  $Y$ , we will prove that

$$E((Y - a - bX)^2) \geq E((Y - \hat{Y})^2)$$

Note the following relation:

$$E((Y - \hat{Y})(\hat{Y} - a - bX)) = 0 \forall a, b$$

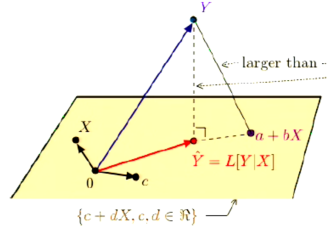
We have shown that error is orthogonal to all linear functions.  $\hat{Y}$  and  $-a - bX$  are both linear functions, so adding them produces a linear function.

Now we use this fact to prove the inequality:

$$\begin{aligned} E((Y - a - bX)^2) &= E((Y - \hat{Y} + \hat{Y} - a - bX)^2) \\ &= E((Y - \hat{Y})^2) + E((\hat{Y} - a - bX)^2) + 0 \\ &\geq E((Y - \hat{Y})^2) \end{aligned} \quad \text{as } E((\hat{Y} - a - bX)^2) \geq 0$$

We have shown that the mean square error from using  $a + bX$  is always going to be larger than if we used  $\hat{Y}$ .

The diagram below offers a visualization of the above facts (from Walrand)



It also demonstrates the **Projection Property**. If the plane is the set of all linear functions of  $X$ , then  $Y - \hat{Y}$  is orthogonal (perpendicular) to any linear function of  $X$  (or the plane).

Here is another derivation of  $\hat{Y}$  using a calculus approach.

## 2.5 Second Proof of $L[Y|X]$

Assume that  $E(X) = E(Y) = 0$ .

$$\hat{Y} = E(Y) + \frac{\text{Cov}(X, Y)}{\text{Var}(X)} \cdot (X - E(X)) = \frac{\text{Cov}(X, Y)}{\text{Var}(X)} \cdot X$$

We want to find  $a$  and  $b$  such that  $g(a, b)$  is minimized, where  $\hat{Y} = a + bX$

$$\begin{aligned} g(a, b) &= E((Y - a - bX)^2) \\ &= E(Y^2 + a^2 + b^2X^2 - 2aX - 2bXY + 2bX) \\ &= a^2 + E(Y^2) + b^2E(X^2) - 2aE(Y) - 2bE(XY) + 2bE(X) \\ &= a^2 + E(Y^2) + b^2E(X^2) - 2bE(XY) \end{aligned} \quad \text{X, Y are zero mean}$$

To find the minimum value of  $g(a, b)$  take the partial derivatives with respect to  $a$  and  $b$  and set them both to zero. Then solve to find  $a$  and  $b$ .

$$0 = \frac{\partial}{\partial a} g(a, b) = 2a \rightarrow a = 0$$

$$0 = \frac{\partial}{\partial b} g(a, b) = 2bE(X^2) - 2E(XY) \rightarrow b = \frac{E(XY)}{E(X^2)} = \frac{\text{Cov}(X, Y)}{\text{Var}(X)}$$

The last line follows from the fact that  $X$  and  $Y$  are zero mean:

$$\text{Cov}(X, Y) = E(XY) - E(X)E(Y) = E(XY) - 0 \cdot 0 = E(XY)$$

$$\text{Var}(X) = E(X^2) - E(X)^2 = E(X^2) - 0 = E(X^2)$$

$$\hat{Y} = a + bX = 0 + \frac{\text{Cov}(X, Y)}{\text{Var}(X)} \cdot X$$

This proves that  $\hat{Y}$  is the best linear estimate for  $Y$ .

In the above proof we forced  $X$  and  $Y$  to be zero mean. We can use a similar proof technique with any  $X$  and  $Y$ .

## 2.6 Third Proof of $L[Y|X]$

---

This applies the second proof to  $X$  and  $Y$  that are not zero mean. This proof is not covered in lecture and basically involves maneuvering terms around, so read it if you're curious.

$$\begin{aligned} Y - a - bX &= Y - E(Y) - (a - E(Y)) - b(X - E(X)) - bE(X) && \text{add some magic terms} \\ &= Y - E(Y) - (a - E(Y) + bE(X)) - b(X - E(X)) && \text{rearrange terms} \\ &= Y - E(Y) - c - b(X - E(X)) && c = a - E(Y) + bE(X) \end{aligned}$$

Now apply results from the Second Proof of  $L[Y|X]$ , since  $Y - E(Y)$  and  $X - E(X)$  are both zero mean.

$$\begin{aligned} c &= 0 \rightarrow 0 = a - E(Y) + bE(X) \rightarrow a = E(Y) - bE(X) \\ b &= \frac{\text{Cov}(X - E(X), Y - E(Y))}{\text{Var}(X - E(X))} = \frac{\text{Cov}(X, Y)}{\text{Var}(X)} \end{aligned}$$

Now plug in the above results into  $a + bX$ :

$$\begin{aligned} a + bX &= E(Y) - bE(X) + bX \\ &= E(Y) - b(-E(X) + x) = E(Y) + \frac{\text{Cov}(X, Y)}{\text{Var}(X)}(X - E(X)) \end{aligned}$$

## 2.7 Some General Comments about LLSE

---

$$L[Y|X] = E(Y) + \frac{\text{Cov}(X, Y)}{\text{Var}(X)} \cdot (X - E(X))$$

1. Covariance controls the slope of LLSE (variance is always positive). If the covariance is positive, then as  $X$  increases, our prediction for  $Y$  will also increase.
2. If you know the distributions of  $X$  and  $Y$ , you can calculate each component of the above equation.
3. If  $X$  and  $Y$  are independent, then the best guess is  $E(Y)$
4. If we have no observations, the best choice for an estimate of  $Y$  is  $E(Y)$ . In this case, the mean error would be:

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

5. Observing  $X$  will reduce the error. The error depends on how  $X$  and  $Y$  are correlated and is defined as follows:

$$E((Y - L[Y|X])^2) = \text{Var}(Y) - \frac{\text{Cov}(X, Y)^2}{\text{Var}(X)}$$

Note that both the numerator and denominator in the fraction are nonnegative. Since the error of using  $E(Y)$  is  $\text{Var}(Y)$ , the error of using  $L[Y|X]$  will always be smaller than if we just used  $E(Y)$ .



## 2.8 Bayesian and non Bayesian

---

We motivated this section with a non Bayesian example. Then we used Bayesian stuff in all of our proofs. Now we explain that Bayesian and non Bayesian are really the same thing.

### non Bayesian

Given samples:  $\{(X_n, Y_n) \mid n = 1, \dots, N\}$  We want

$$\hat{Y} = a + bX$$

where  $a, b$  minimize

$$\sum_{n=1}^N (Y_n - a - bX_n)^2$$

In other words, we minimize the square of the distance between the line and the points. We take the square of the distance because taking the absolute value of the difference is more complicated.

Note that in the above formulation, we made no assumptions about the distribution of  $X$ .

### Bayesian

Given  $X, Y$  with known distribution. We want:

$$\hat{Y} = a + bX = L[Y|X]$$

where

$$g(a, b) = E((Y - a - bX)^2)$$

is minimized.

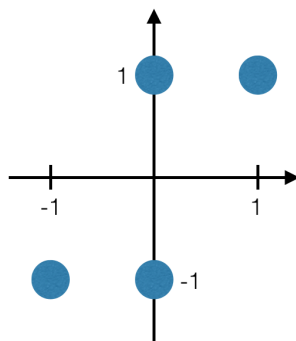
We can turn non Bayesian into Bayesian as follows:

$$\frac{1}{N} \sum (Y_n - a - bX_n)^2 = E((Y - a - bX)^2)$$

$$(X, Y) = (X_n, Y_n) \text{ w. p. } \frac{1}{N} \text{ for } n = 1, \dots, N$$

### 3 Example of Finding LLSE (non Bayesian)

Find the linear least square estimate of the following joint distribution.



We turn the problem from non Bayesian into Bayesian by assuming that the 4 points are equally likely values of  $(X, Y)$ .

$$E(X) = -1 \cdot \frac{1}{4} + 0 \cdot \frac{1}{2} + 1 \cdot \frac{1}{4}$$

$$E(Y) = -1 \cdot \frac{1}{2} + 1 \cdot \frac{1}{2} = 0$$

$$E(X^2) = (-1)^2 \cdot \frac{1}{4} + (0)^2 \cdot \frac{1}{2} + (1)^2 \cdot \frac{1}{4} = \frac{1}{2}$$

$$E(XY) = 1 \cdot \frac{1}{2} = \frac{1}{2}$$

$$\text{Var}(X) = E(X^2) - E(X)^2 = \frac{1}{2} - 0 = \frac{1}{2}$$

$$\text{Cov}(X, Y) = E(XY) - E(X)E(Y) = \frac{1}{2} - 0 = \frac{1}{2}$$

Therefore we get:

$$L[Y|X] = 0 + \frac{\frac{1}{2}}{\frac{1}{2}} \cdot X = X$$

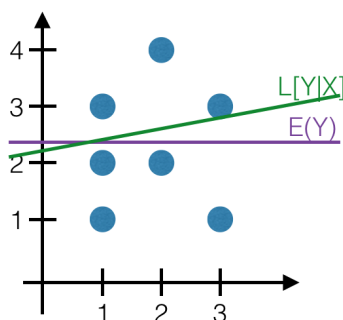
## 4 Conditional Expectation

### Motivation

There are situations where the relationship between  $X$  and  $Y$  is non linear. Our goal is to find  $g(\cdot)$  such that  $g(X)$  is the best guess about  $Y$  given  $X$  where  $g(\cdot)$  is not restricted to a line.

### 4.1 Motivating Example

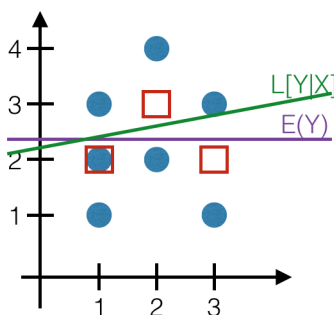
Examine the following joint distribution of  $X$  and  $Y$ .



Without any observations, our best guess for  $Y$  is  $E(Y) = 2.3$ . Assuming we observe  $X$ , we can construct the following linear estimate:

$$L[Y|X] = a + bX = 2.1 + 0.1X$$

What if we know that  $X$  is 1? What would be our best guess for  $Y$  in that case? 2! Since there are 3 values that  $Y$  can take on one we know that  $X$  is 1, we find the average of those values. So once we narrow down  $X$ , we examine a new distribution. Repeating this process for each value of  $X$  gives us the following diagram:



Compare the each of the new expected values with what the linear equation gives us. Which is closer?

## 4.2 Defining Conditional Expectation

---

If  $X$  and  $Y$  are random variables in  $\Omega$  then the **conditional expectation** is defined as follows:

$$\begin{aligned} E(Y|X) &= g(X) \\ \text{where, } g(X) &= E(Y|X = x) \\ &= \sum_y y \cdot P[Y = y|X = x] \\ \text{with } P[Y = y|X = x] &= \frac{P[X = x, Y = y]}{P[X = x]} \end{aligned}$$

### Theorem

$E(Y|X)$  is the best guess about  $Y$  given  $X$ .

$$E((Y - g(X))^2) \geq E((Y - E(Y|X))^2)$$

Using any other function to predict  $Y$  will give a larger error.

## 4.3 Projection Property Revisited

---

**Projection Property:**  $E((Y - E(Y|X))f(X)) = 0 \forall f(\cdot)$

This is equivalent to the following,

$$\begin{aligned} E((Y - E(Y|X))f(X)) &= E(Y \cdot f(X)) - E(E(Y|X) \cdot f(X)) \\ &\rightarrow \\ E(Y \cdot f(X)) &= E(E(Y|X)f(X)) \end{aligned}$$

### Proof:

$$\begin{aligned} E(E(Y|X) \cdot f(X)) &= \sum_x E(Y|X = x) \cdot f(x) \cdot P[X = x] \\ &= \sum_x \left( \sum_y y \cdot f(x) \cdot P[Y = y|X = x] \right) \cdot P[X = x] \\ &= \sum_x \sum_y y \cdot f(x) P[X = x, Y = y] \quad \text{using } \frac{P[X = x, Y = y]}{P[X = x]} \\ &= E(Y \cdot f(X)) \end{aligned}$$

Therefore, if we have the conditional expectation, we can take its mean value and get the mean value of  $Y$ .

**4.4 CE = MMSE**

---

Finally we prove the fact that  $E(Y|X)$  is the best guess about  $Y$  based on  $X$ . We use a technique similar to before.

Recall,

$$\begin{aligned} E(Y|X) &= g(X) \\ \text{where, } g(X) &= E(Y|X = x) \\ &= \sum_y y \cdot P[Y = y|X = x] \\ \text{with } P[Y = y|X = x] &= \frac{P[X = x, Y = y]}{P[X = x]} \end{aligned}$$

We want to show that  $g(X)$  minimizes  $E((Y - g(X))^2)$ . Let  $h(X)$  be any other function.

$$\begin{aligned} E((Y - h(X))^2) &= E((Y - g(X) + g(X) - h(X))^2) && \text{add in some magic terms} \\ &= E((Y - g(X))^2) + E((g(X) - h(X))^2) \\ &\quad + 2E((Y - g(X))(g(X) - h(X))) && \text{linearity of expectation} \end{aligned}$$

We know that  $E((Y - g(X))(g(X) - h(X))) = 0$  by the Projection Property and  $E((g(X) - h(X))^2) \geq 0$ .

Therefore,

$$E((Y - h(X))^2) \geq E((Y - g(X))^2)$$