

# Lecture 8: Conditional Expectation

Ziyu Shao

School of Information Science and Technology  
ShanghaiTech University

December 10, 2024

# Outline

- 1 Conditional Expectation: Given An Event
- 2 Conditional Expectation: Given A Random Variable
- 3 Prediction & Estimation
- 4 Application Case: Kalman Filter

# Outline

- 1 Conditional Expectation: Given An Event
- 2 Conditional Expectation: Given A Random Variable
- 3 Prediction & Estimation
- 4 Application Case: Kalman Filter

# Conditional PMF

- Let  $A$  be an event with positive probability. If  $X$  is a discrete r.v., then the conditional PMF of  $X$  given  $A$  is

$$P_{X|A}(x) = P(X = x|A) = \frac{P(\{X = x\} \cap A)}{P(A)}.$$

- Bayes' Rule:

$$P_{X|A}(x) = P(X = x|A) = \frac{P(A|X = x)P(X = x)}{P(A)}.$$

- LOTP: with a partition  $A_1, \dots, A_n$ , each  $A_i$  with a positive probability  $P(A_i) > 0$ ,  $i = 1, 2, \dots, n$ :

$$P(X = x) = \sum_{i=1}^n P_{X|A_i}(x)P(A_i).$$



# Conditional PDF

- Let  $A$  be an event with positive probability. If  $X$  is a continuous r.v., then the *conditional PDF of  $X$  given  $A$*  is

$$f_{X|A}(x) = (P(X \leq x|A))'.$$

- LOTP: with a partition  $A_1, \dots, A_n$ , each  $A_i$  with a positive probability  $P(A_i) > 0$ ,  $i = 1, 2, \dots, n$ :

$$f_X(x) = \sum_{i=1}^n P(A_i) f_{X|A_i}(x).$$

# Conditional PDF

- Bayes' Rule: given an event  $A$  with  $P(A) > 0$ , then

$$f_{X|A}(x) = \frac{P(A|X=x)}{P(A)} \cdot f_X(x).$$

- Bayes' Rule: given event  $A = "a \leq X \leq b"$  and  $P(A) > 0$ , then

$$\begin{aligned} f_{X|A}(x) &= \frac{\mathbf{1}_{x \in [a,b]}}{P(A)} \cdot f_X(x) \\ &= \begin{cases} \frac{f_X(x)}{P(A)} & \text{if } a \leq x \leq b \\ 0 & \text{otherwise} \end{cases} \end{aligned}$$

# Conditional Expectation Given An Event

## Definition

Let  $A$  be an event with positive probability. If  $Y$  is a discrete r.v., then the *conditional expectation of  $Y$  given  $A$*  is

$$\underline{E(Y|A)} = \sum_y \underline{y \cdot P(Y = y|A)} = \sum_y y \cdot P_{Y|A}(y),$$

where the sum is over the support of  $Y$ . If  $Y$  is a continuous r.v. with PDF  $f$ , then

$$E(Y|A) = \int_{-\infty}^{\infty} y \cdot \underline{f_{Y|A}(y)} dy.$$

# LOTUS Given An Event

## Definition

Let  $A$  be an event with positive probability and  $g$  is a function from  $\mathbf{R}$  to  $\mathbf{R}$ . If  $Y$  is a discrete r.v., then the *conditional expectation of  $g(Y)$  given  $A$*  is

$$\underline{E(g(Y)|A)} = \sum_y \underline{g(y) \cdot P_{Y|A}(y)},$$

where the sum is over the support of  $Y$ .

If  $Y$  is a continuous r.v. with PDF  $f$ , then

$$E(g(Y)|A) = \int_{-\infty}^{\infty} \underline{g(y) \cdot f_{Y|A}(y)} dy.$$

# Example

Method 1: 1<sup>o</sup>. event  $A = "X > 1"$ ,  $P(A) = P(X > 1) = \int_1^{\infty} \lambda e^{-\lambda x} dx = e^{-\lambda}$

$$f_{X|A}(x) = \frac{f(x)}{P(A)} = \lambda e^{-\lambda(x-1)}, \quad x > 1.$$

$$2^o. \quad \underline{E[X|X > 1]} = \int_1^{\infty} x \cdot f_{X|A}(x) dx = \int_1^{\infty} x \cdot \lambda e^{-\lambda(x-1)} dx = 1 + \frac{1}{\lambda}$$

$$\underline{E[X^2|X > 1]} = \int_1^{\infty} x^2 \cdot f_{X|A}(x) dx = \int_1^{\infty} x^2 \cdot \lambda e^{-\lambda(x-1)} dx = \frac{1 + 2/\lambda + 1/\lambda^2}{\lambda^2}$$

$$\Rightarrow \underline{\text{Var}(X|X > 1)} = E[X^2|X > 1] - (E[X|X > 1])^2 = \frac{1}{\lambda^2}$$

Let  $X \sim \text{Expo}(\lambda)$ , find  $E(X|X > 1)$  and  $\text{Var}(X|X > 1)$ .

$$E(X) = \frac{1}{\lambda}, \quad \text{Var}(X) = \frac{1}{\lambda^2};$$

Method 2:

$$\forall s, t > 0, \quad P(X > s+t | X > t) = P(X > s) \quad ; \quad t=1$$

$$\Rightarrow P(X > s+1 | X > 1) = P(X > s) \Rightarrow \underline{P(X-1 > s | X > 1) = P(X > s)}$$

$$1^o \quad E[X|X > 1] = E[X-1+1|X > 1] = 1 + E[X-1|X > 1]$$

$$= 1 + E[X] = 1 + \frac{1}{\lambda}$$

$$2^o. \quad \text{Var}(X|X > 1) = \text{Var}(X-1|X > 1)$$

$$= \text{Var}(X) = \frac{1}{\lambda^2}$$

$$\underline{X-1 | X > 1 \wedge X}$$

$$E(X-1|X > 1) = E(X)$$

$$\underline{\text{Var}(X-1|X > 1) = \text{Var}(X)}$$

# Solution

# Motivation of Conditional Expectation

- Conditional expectation is a powerful tool for calculating expectations: first-step analysis
- Conditional expectation allows us to predict or estimate unknowns based on whatever evidence is currently available.
- Conditional Expectation given an event:  $E(Y|A)$
- Conditional Expectation given a random variable:  $E(Y|X)$

# Life Expectancy

$T$  : Life span.

$$E(T) = 70;$$

---

$$E[T | T \geq 20] \neq E(T)$$

$$T \wedge \text{Exp}(\cdot) \Rightarrow 20 + E(T)$$



# Law of Total Expectation

LOTE

W.L.O.G.

Discrete.

$$E(Y|A_i) = \sum y P(Y=y|A_i)$$

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

$$= \sum_{y=y}^n y \cdot \sum P(Y=y|A_i) \cdot P(A_i)$$

$$= \sum y \cdot \sum_{i=1}^n P(Y=y, A_i)$$

LOTE  $\rightarrow$  LOTP:  $Y = I_B$

$$\Rightarrow P(B) = E(I_B) = E(Y) \stackrel{\text{LOTE}}{=} \sum_{i=1}^n E(Y|A_i) \cdot P(A_i)$$

Theorem

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

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

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

$$= \sum y \cdot P(Y=y) = E(Y)$$

Memoryless

$X \sim \text{Exp}(\lambda)$ ; And  $E[X|X \leq 1]$ ?

$X-1|X > 1 \sim X$

$$E[X] = E[X|X > 1] \cdot P(X > 1) + E[X|X \leq 1] \cdot P(X \leq 1)$$

$$\frac{1}{\lambda} = [1 + \frac{1}{\lambda}] \cdot e^{-\lambda} + E[X|X \leq 1] \cdot (1 - e^{-\lambda})$$

# Geometric Expectation Redux

$X \sim \text{Geom}(p)$ , Find  $E[X]$ .

① First step Analysis. (Conditioning on the outcome of the first toss.

$$O_1 = H \text{ or } T.$$

$$\begin{aligned} \textcircled{2} \quad E(X) &\stackrel{\text{LOTE}}{=} \underbrace{E(X|O_1=H)} \cdot \underbrace{P(O_1=H)} + \underbrace{E(X|O_1=T)} \cdot \underbrace{P(O_1=T)} \\ &= 0 \cdot p + \underbrace{(1 + E(X)) \cdot (1-p)} \end{aligned}$$

$$\Rightarrow E(X) = \frac{1-p}{p}$$

Coin Tosses

# of tosses before the first Heads appearance.

H: Head ; T: Tail.

# Outline

- 1 Conditional Expectation: Given An Event
- 2 Conditional Expectation: Given A Random Variable
- 3 Prediction & Estimation
- 4 Application Case: Kalman Filter

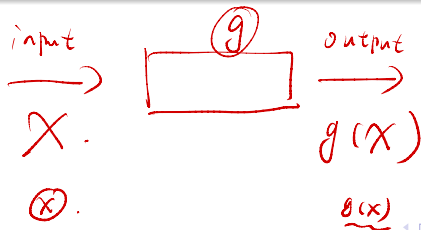
# Conditional Expectation Given An R.V.

$$g(x) = E(Y | X=x) \quad \text{Real Number} \quad \underline{\text{estimate.}}$$

$$g(X) = E(Y | X) \quad \text{R.V.} \quad \text{estimator.}$$

## Definition

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



# Remark

- $E(Y|X)$  is a function of  $X$ , and it is a random variable.
- It makes sense to compute  $E(E(Y|X))$  and  $Var(E(Y|X))$ .

## Example: Stick Length



$$\textcircled{1} X \sim \text{unif}(0,1): Y|X=x \sim \text{unif}(0,x) \Rightarrow E(Y|X=x) = \frac{x}{2} \\ \Rightarrow E(Y|X) = g(X) = \frac{X}{2} = g(x).$$

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

Suppose we have a stick of length 1 and break the stick at a point X chosen uniformly at random. Given that  $X = x$  we then choose another breakpoint Y uniformly on the interval  $[0, x]$ . Find  $E(Y|X)$ , and its mean and variance.

$$\text{Var}[E(Y|X)] = \text{Var}\left(\frac{X}{2}\right) = \frac{1}{4} \text{Var}(X) = \frac{1}{4} \cdot \frac{1}{12} = \frac{1}{48}$$

Find  $E(Y|X)$

$$\Rightarrow \textcircled{1} g(x) = E(Y|X=x).$$

$$\textcircled{2} g(x) \rightarrow g(X) = E(Y|X) \\ x \rightarrow X.$$

# Solution

# Dropping What's Independent

## Theorem

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

$$g(x) = E[Y|X=x] = E[Y], \quad \forall x.$$

$$\Rightarrow g(x) = \underline{E(Y)}$$

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



# Taking Out What's Known

$$\begin{aligned} 1^\circ. g(x) &= E[h(X)Y | X=x] = E[\underbrace{h(x)} Y | X=x] \\ &= \underbrace{h(x)} \cdot E[Y | X=x] \end{aligned}$$

$$2^\circ. g(X) = \underbrace{h(X)} \cdot \underbrace{E[Y|X]}$$

## Theorem

For any function  $h$ ,

$$E(\underbrace{h(X)} Y | \underbrace{X}) = \underbrace{h(X)} \underbrace{E(Y|X)}$$

# Linearity

$$g(x) = E[Y_1 + Y_2 | X=x] = E[Y_1 | X=x] + E[Y_2 | X=x]$$

$$\Rightarrow g(X) = E[Y_1 | X] + E[Y_2 | X]$$

## Theorem

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

## Example

1°. By iid.  $E[X_1|S_n] = E[X_2|S_n] = \dots = E[X_n|S_n]$

2°. By Linearity.  $E[X_1|S_n] + E[X_2|S_n] + \dots + E[X_n|S_n]$

$$g(s) = E[S_n|S_n=s] = s$$

$$g(S_n) = S_n$$

$$= E[X_1 + X_2 + \dots + X_n | S_n]$$

$$= E[S_n | S_n] = S_n$$

Let  $X_1, \dots, X_n$  be i.i.d., and  $S_n = X_1 + \dots + X_n$ . Find  $E(X_1|S_n)$ .

3°.  $n E[X_1|S_n] = S_n$

$$\Rightarrow E[X_1|S_n] = \frac{1}{n} S_n$$

# Adam's Law

The Law of Iterated Expectation

The Tower Rule.

The Smoothing Theorem

## Theorem

W.L.O.G.  $X$  and  $Y$  are both discrete r.v.s.

For any r.v.s  $X$  and  $Y$ , 1<sup>o</sup>.  $g(X) = E(Y|X)$ ;  $g(x) = E(Y|X=x)$

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

2<sup>o</sup>. LHS.  $E[E(Y|X)] = E[g(X)] = \sum_x g(x) \cdot P(X=x)$

$= \sum_x \left( \sum_y y \cdot P(Y=y|X=x) \right) \cdot P(X=x) = \sum_y y \cdot \left[ \sum_x P(Y=y|X=x) \cdot P(X=x) \right]$

$= \sum_y y \cdot P(Y=y) = E(Y)$  RHS.

# Proof

# Adam's Law with Extra Conditioning

$$\hat{P}(\cdot) = P(\cdot|Z) \quad ; \quad \hat{E} = E(\cdot|Z)$$

Adam's Law :  $E[E(Y|X)] = E(Y)$

$$\hat{E}[\hat{E}(Y|X)] = \hat{E}(Y)$$

## Theorem

For any r.v.s  $X, Y, Z$ , we have

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

$$\Downarrow \\ \underline{E(Y|Z)}$$

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

$$\hat{E}(Y|X)$$

$$= \underline{E(Y|X, Z)}$$

Conditional Variance 1<sup>o</sup>.  $\text{Var}(Y) = \underline{E}[(Y - \underline{E}(Y))^2]$

$$\hat{E}(\cdot) = E(\cdot|X)$$

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

## Definition

The conditional variance of  $Y$  given  $X$  is

$$\underline{\text{Var}}(Y|X) = \underline{E}\left(\left(Y - \underline{E}(Y|X)\right)^2 \mid X\right).$$

This is equivalent to

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

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

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

# Eve's law

## Theorem

For any r.v.s  $X$  and  $Y$ ,

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

The ordering of  $E$ 's and  $\text{Var}$ 's on the right-hand side spells EVVE, whence the name Eve's law. Eve's law is also known as the law of total variance or the variance decomposition formula.



# Proof

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

$$\textcircled{1} \quad \underline{g(X) = E(Y|X)} \quad \Rightarrow \text{Adam's Law} \quad \underline{E(g(X)) = E[E(Y|X)] = E(Y)}$$

$$\begin{aligned} \textcircled{2} \quad \underline{E[\text{Var}(Y|X)]} &= E[E(Y^2|X) - \underbrace{(E(Y|X))^2}_{g^2(x)}] = E[E(Y^2|X) - g^2(x)] \\ &= \underline{E[E(Y^2|X)]} - E[g^2(x)] \stackrel{\text{Adam's Law}}{=} \underline{E(Y^2)} - \underline{E(g^2(x))} \end{aligned}$$

$$\begin{aligned} \textcircled{3} \quad \underline{\text{Var}[E(Y|X)]} &= \text{Var}(g(X)) = E[g^2(x)] - \underline{E^2(g(x))} \\ &= \underline{E[g^2(x)]} - \underline{E^2(Y)} \end{aligned}$$

$$\begin{aligned} \textcircled{2} + \textcircled{3} \quad \underline{E[\text{Var}(Y|X)] + \text{Var}[E(Y|X)]} &= \underline{E(Y^2) - E^2(Y)} \\ &= \underline{\text{Var}(Y)} \end{aligned}$$

# Proof

## Example: Random Sum

$N, X_j$  independ.

$$\textcircled{1} E(X) : \quad E(X|N=n) = E\left[\sum_{j=1}^N X_j | N=n\right] = E\left[\sum_{j=1}^n X_j | N=n\right]$$

$$\Rightarrow \underline{E(X|N)} = g(N) = N \cdot \mu. \quad = E\left[\sum_{j=1}^N X_j\right] = \sum_{j=1}^N E(X_j) = \underline{n \cdot \mu} = g(n)$$

$$\Rightarrow E(X) = E[E(X|N)] = E[N \cdot \mu] = \mu \cdot E(N).$$

A store receives  $N$  customers in a day, where  $N$  is an r.v. with finite mean and variance. Let  $X_j$  be the amount spent by the  $j$ th customer at the store. Assume that each  $X_j$  has mean  $\mu$  and variance  $\sigma^2$ , and that  $N$  and all the  $X_j$  are independent of one another. Find the mean and variance of the random sum  $X = \sum_{j=1}^N X_j$ , which is the store's total revenue in a day, in terms of  $\mu$ ,  $\sigma^2$ ,  $E(N)$ , and  $\text{Var}(N)$ .

$$\textcircled{2} \text{Var}(X|N=n) = \text{Var}\left(\sum_{j=1}^N X_j | N=n\right) = \text{Var}\left(\sum_{j=1}^n X_j | N=n\right)$$
$$= \text{Var}\left(\sum_{j=1}^n X_j\right) = \sum_{j=1}^n \text{Var}(X_j) = \underline{n \sigma^2}$$

$$\Rightarrow \text{Var}(X|N) = N \cdot \sigma^2.$$

## Solution

③ By Eve's Law

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

---

# Solution

# Outline

- 1 Conditional Expectation: Given An Event
- 2 Conditional Expectation: Given A Random Variable
- 3 Prediction & Estimation**
- 4 Application Case: Kalman Filter

# Basic Problem

$$\hat{Y} = g(X)$$

- Estimate  $Y$  from the observed value  $X$
- Choose the estimator (inference function)  $g(\cdot)$  to minimize the expected error  $E(c(Y, g(X)))$
- $c(Y, \hat{Y})$  is the cost of guessing  $\hat{Y}$  when the actual value is  $Y$ .
- When  $c(Y, \hat{Y}) = \|Y - \hat{Y}\|^2$ , the best guess is called “the least square estimate (LSE)” estimate of  $Y$  given  $X$ .
- Further, if the function  $g(\cdot)$  is restricted to be linear, i.e., of the form  $a + bX$ , it is called “the Linear Least Square Estimate (LLSE)” estimate of  $Y$  given  $X$ .
- Further, if the function  $g(\cdot)$  can be arbitrary, it is called “the Minimum Mean Square Estimate (MMSE)” estimate of  $Y$  given  $X$ .

# Linear Least Square Estimate

$$\textcircled{1} \quad f(a,b) = E[(Y - a - bX)^2] = a^2 + E(Y^2) + b^2 E(X^2) - 2aE(Y) + 2abE(X) - 2bE(XY)$$

$$\textcircled{1} \quad \frac{\partial f(a,b)}{\partial a} = 2a - 2E(Y) + 2bE(X) = 0 \Rightarrow \underline{a + bE(X) = E(Y)}$$

## Theorem

The Linear Least Square Estimate (LLSE) of  $Y$  given  $X$ , denoted by  $L[Y|X]$ , is the linear function  $a + bX$  that minimizes  $E[(Y - a - bX)^2]$ . In fact,

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

$$\textcircled{2} \quad \frac{\partial f(a,b)}{\partial b} = 2bE(X^2) + 2aE(X) - 2E(XY) \Rightarrow \underline{aE(X) + bE(X^2) = E(XY)}$$



Proof  $\Rightarrow b = \frac{\text{Cov}(X, Y)}{\text{Var}(X)}$ ,  $a = E[Y] - \frac{\text{Cov}(X, Y)}{\text{Var}(X)} \cdot E[X]$

$$\Rightarrow \underline{E[Y|X]} = a + bX = \underline{E[Y]} + \frac{\text{Cov}(X, Y)}{\text{Var}(X)} [X - E[X]]$$

Hessian Matrix:  $H = \begin{bmatrix} \frac{\partial^2 f(a, b)}{\partial a^2} & \frac{\partial^2 f(a, b)}{\partial a \partial b} \\ \frac{\partial^2 f(a, b)}{\partial a \partial b} & \frac{\partial^2 f(a, b)}{\partial b^2} \end{bmatrix} \succcurlyeq 0$

$$H \succcurlyeq 0 \Leftrightarrow \underline{z^T H z} \geq 0$$

$$H = \begin{bmatrix} 2 & 2E[X] \\ 2E[X] & 2E[X^2] \end{bmatrix} = 2 \begin{bmatrix} 1 & E[X] \\ E[X] & E[X^2] \end{bmatrix}, \quad z = (z_1, z_2)^T$$

$$z^T H z = (z_1, z_2) H \begin{pmatrix} z_1 \\ z_2 \end{pmatrix} = 2 [z_1^2 + 2z_1 z_2 E[X] + z_2^2 E[X^2]]$$

$$= 2 [ (z_1 + z_2 E[X])^2 + z_2^2 (E[X^2] - E^2[X]) ] = 2 [ (z_1 + z_2 E[X])^2 + z_2^2 \text{Var}(X) ] \geq 0$$

# Proof

# Data-Driven LLSE: Linear Regression

$L_k(Y|X)$

- $L[Y|X] = E(Y) + \frac{\text{Cov}(X,Y)}{\text{Var}(X)}(X - E(X))$
- Now we only have  $k$  i.i.d samples:  $(X_1, Y_1), \dots, (X_k, Y_k)$
- Use sample mean to replace expectation

$$E(X) \leftarrow E_k(X) = \frac{1}{k} \sum_{j=1}^k X_k$$

$E \leftarrow$  Sample mean.

$k \rightarrow \infty$

$$E(Y) \leftarrow E_k(Y) = \frac{1}{k} \sum_{j=1}^k Y_k$$

SSLN ;

LR  $\rightarrow$  LLSE

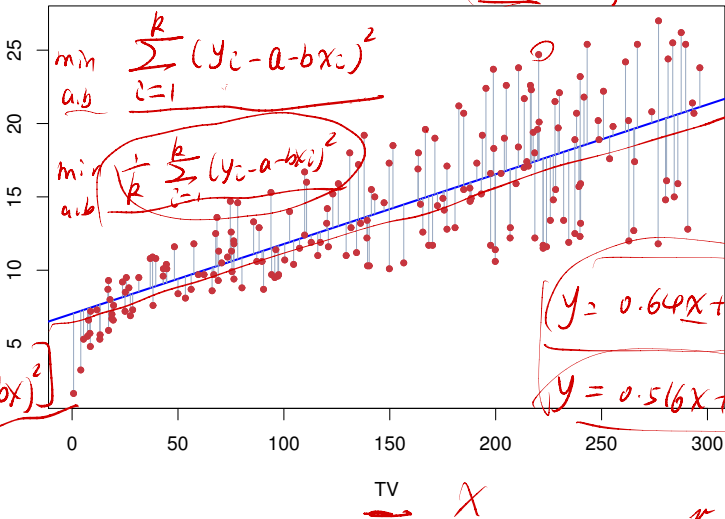
$$\text{Cov}(X, Y) = E(XY) - E(X)E(Y) \leftarrow \frac{1}{k} \sum_{j=1}^k X_k Y_k - E_k(X)E_k(Y)$$

$$\text{Var}(X) = E(X^2) - (E(X))^2 \leftarrow \frac{1}{k} \sum_{j=1}^k X_k^2 - (E_k(X))^2$$

# Data-Driven LLSE: Linear Regression

$$y = 0.512x + 0.62$$

$(X_i, y_i)$  data  
Son's height  
Father's height



# Minimum Mean Square Error Estimator

$$\text{LLSE: } \min_{\hat{Y}} E[(Y - \hat{Y})^2], \hat{Y} = a + bX \Rightarrow \hat{Y}^* = \underline{L(Y|X)}$$

$$\text{MMSE: } \min_{\hat{Y}} E[(Y - \hat{Y})^2], \hat{Y} = g(X) \Rightarrow \hat{Y}^* = E[Y|X]$$

## Theorem

The MMSE of  $Y$  given  $X$  is given by

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

$E(X^2) < \infty$

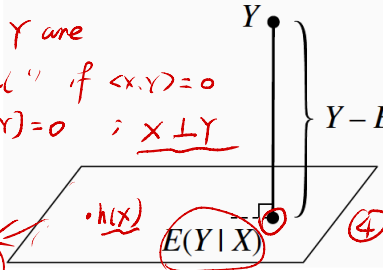
# Geometric Perspective of Conditional Expectation

① inner product  $\langle X, Y \rangle = E(X \cdot Y)$  ,  $\|X\| = \sqrt{\langle X, X \rangle} = \sqrt{E(X^2)}$

cos  $\theta = \frac{\langle X, Y \rangle}{\|X\| \cdot \|Y\|}$  ;  $\text{dist}(X, Y) = \sqrt{\langle X - Y, X - Y \rangle} = \sqrt{E((X - Y)^2)}$

② X and Y are "orthogonal" if  $\langle X, Y \rangle = 0$   
 $\Leftrightarrow E(X \cdot Y) = 0$  ;  $X \perp Y$

③  $\text{COV}(X, Y) = E(X \cdot Y) - E(X) \cdot E(Y)$   
if  $E(X) = 0$  or  $E(Y) = 0$  or both.  
 $\Rightarrow \text{COV}(X, Y) = E(X \cdot Y)$   
 $\Rightarrow$  Uncorrelated  $\Leftrightarrow$  Orthogonal



space of  $h(X)$  ( $\cup h(\cdot)$ )

④  $Y - E(Y|X) \perp h(X)$   
 $\cup h(\cdot)$

⑤  $E(Y|X)$ : a projection of Y onto the space of  $(h(X) \cup h(\cdot))$

⑥  $L(Y|X)$ : a projection of Y onto the  $L(X) = \{a + bX, a, b \in \mathbb{R}\}$   
LLSB

Projection Interpretation

$$\begin{aligned}
 1^\circ. \quad E[\underbrace{Y - E(Y|X)}] &= E[Y] - \underbrace{E[E(Y|X)]} \\
 &= E[Y] - E[Y] \\
 &= 0
 \end{aligned}$$

$$\text{Cov}(Y - E(Y|X), h(X)) = E[(Y - E(Y|X))h(X)]$$

## Theorem

For any function  $h$ , the r.v.  $Y - E(Y|X)$  is uncorrelated with  $h(X)$ .  
 Equivalently,

$$E((Y - E(Y|X))h(X)) = 0.$$

(This is equivalent since  $E(Y - E(Y|X)) = 0$ , by linearity and Adam's law.)

2<sup>o</sup>. Show  $Y - E(Y|X) \perp h(X)$

Uncorrelated  $\Leftrightarrow$  Orthogonal.

# Proof

$$Y - E(Y|X) \perp h(X)$$

$$\Leftrightarrow \underline{E[(Y - E(Y|X)) \cdot h(X)] = 0}$$

$$= E[Y h(X) - h(X) \cdot E(Y|X)]$$

$$= E[Y h(X)] - \underline{E[h(X) \cdot E(Y|X)]}$$

$$= \underline{E[Y h(X)]} - \underline{E[E[h(X) \cdot Y | X]]}$$

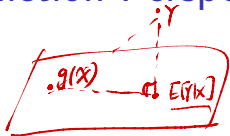
$$= E[Y \cdot h(X)] - E[h(X) \cdot Y] \quad \downarrow \text{Adam's Law}$$

$$= 0$$



# Proof

# Prediction Perspective



$$\underline{E[(Y - E(Y|X))^2]} \leq \underline{E[(Y - \hat{Y})^2]}$$

$$\underline{g(x)}$$

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

- Predict or estimate the future observations or unknown parameters based on data
- $E(Y|X)$  is our **best predictor** of  $Y$  based on  $X$ .
- Best means it is the function of  $X$  with the lowest mean squared error (expected squared difference between  $Y$  and prediction of  $Y$ ).
- It is called the Minimum Mean Square Estimate (MMSE)

Proof 1°.  $\hat{Y}$ : estimator of  $Y$  based on  $X$  ( $\hat{Y} = g(X)$ )

$$E[(Y - \hat{Y})^2] = E[(Y - g(X))^2] ; \quad Y - g(X) = \underbrace{Y - E(Y|X)}_A + \underbrace{E(Y|X) - g(X)}_B$$

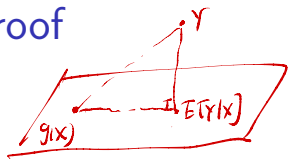
$$2^\circ. E[(Y - g(X))^2] = E[(A+B)^2] = E[A^2] + 2E[AB] + E[B^2]$$
$$= E[(Y - E(Y|X))^2] + E[(E(Y|X) - g(X))^2] + \dots$$

$$2 E[(Y - E(Y|X)) \cdot (E(Y|X) - g(X))]$$

$$3^\circ. h(X) \triangleq \underbrace{E(Y|X) - g(X)} ; \quad \underbrace{Y - E(Y|X)} \perp h(X) \Rightarrow \underbrace{E[(Y - E(Y|X)) \cdot h(X)] = 0}$$

$$\Rightarrow \underbrace{E[(Y - g(X))^2]} = \underbrace{E[(Y - E(Y|X))^2]} + \underbrace{E[(E(Y|X) - g(X))^2]}$$

# Proof



$$\begin{aligned} \text{dist}^2(Y, E(Y|X)) + \text{dist}^2(g(X), E(Y|X)) \\ = \text{dist}^2(Y, g(X)) \end{aligned}$$

---

$$\begin{aligned} 5^\circ. \quad E[(Y - \hat{Y})^2] &= E[(Y - g(X))^2] \\ &= \underbrace{E[(Y - E(Y|X))^2]} + \underbrace{E[(E(Y|X) - g(X))^2]} \\ &\geq \underbrace{E[(Y - E(Y|X))^2]} \quad (\geq 0) \end{aligned}$$

$$\hat{Y}^* = g^*(X) = \underbrace{E(Y|X)}$$

MSE

# Proof

# MMSE for Jointly Normal Random Variables

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

$$\rho = \frac{\text{Cov}(X,Y)}{\sqrt{\text{Var}(X)\text{Var}(Y)}} = \frac{\text{Cov}(X,Y)}{\sigma_X \sigma_Y}$$

## Theorem

Let  $X, Y$  be jointly Normal random variables. Then

$$\underbrace{E[Y|X]}_{\text{MMSE}} = \underbrace{L[Y|X]}_{\text{LLSE}} = E(Y) + \frac{\text{Cov}(X, Y)}{\text{Var}(X)} (X - E(X)).$$

$$\Rightarrow E[Y|X] - E[Y] = \frac{\rho \sigma_X \sigma_Y}{\sigma_X^2} (X - E(X))$$

$$\Rightarrow \frac{E[Y|X] - E[Y]}{\sigma_Y} = \rho \cdot \frac{X - E(X)}{\sigma_X} \Rightarrow$$

$X$ : parent's height.

$Y$ : child's height.

$$\underline{X - E(X) = t \sigma_X}$$

# Remark: Statistical Learning Perspective

$X$ : parent's height

$Y$ : child's height.

$$p > 0$$

$$0 < p < 1$$

$$X \sim E(X) = t \sigma_X$$

$$\Rightarrow \frac{E(Y|X) - E(Y)}{\sigma_Y} = \underbrace{(pt)}_{(t)} < (t) \Rightarrow \underbrace{(E(Y|X) - E(Y))}_{(t \sigma_Y)} < (t \sigma_Y)$$

- In general, MMSE is a highly nonlinear function.
- Adoption of various approximation methods leads to various learning methods

- ▶ Linear regression
- ▶ Logistic regression
- ▶ Polynomial regression
- ▶ Regression with Spline functions
- ▶ Neural network

$$g(x) = a + bx$$

$$g(x) = \frac{1}{1 + e^{-(a+bx)}}$$

$$g(x) = a + bx + cx^2 + dx^3 + \dots + fx^6$$

piecewise polynomial

$$\frac{\log \frac{g(x)}{1-g(x)}}{g(x)} = a + bx$$

logit

regression towards Mean

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

regression function

# Outline

- 1 Conditional Expectation: Given An Event
- 2 Conditional Expectation: Given A Random Variable
- 3 Prediction & Estimation
- 4 Application Case: Kalman Filter



# Milestones in Statistics & Signal Processing

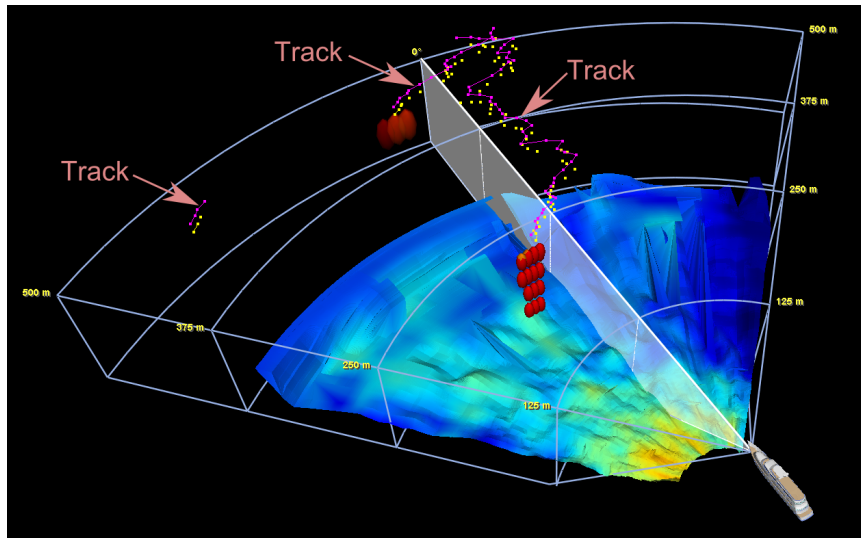
- 1960: Rudolph Emil Kalman (1930-2016) introduced what is known as Kalman filter.



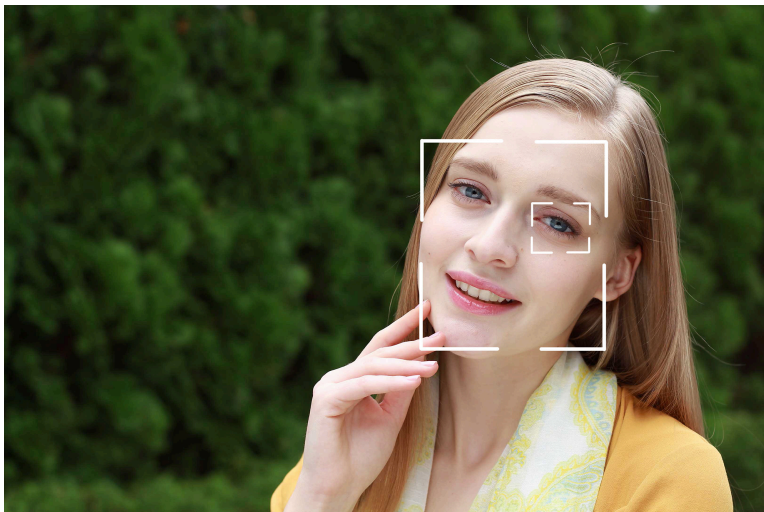
# Widely Applications: Location & Navigation & Map Building



# Widely Applications: Radar Tracking



# Widely Applications: Human Face & Eye Detection Autofocus



# Widely Applications: Animal Eye Detection Autofocus



# Essence of Kalman Filter

Online LLSE

(MMSE)

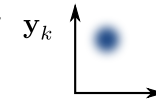
$n$ : index of time.

$M_n = \frac{1}{n} (x_1 + \dots + x_n)$

offline

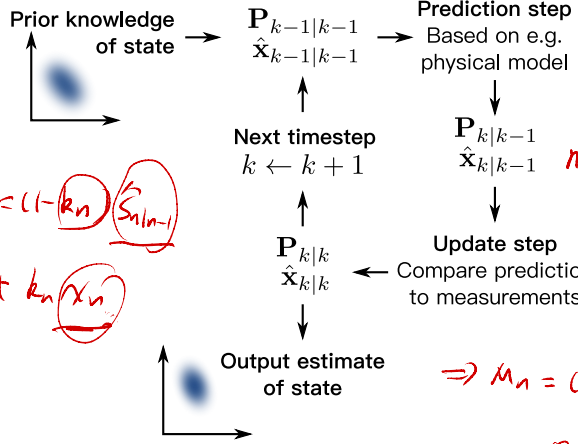
$nM_n = (x_1 + \dots + x_{n-1}) + x_n$

$= (n-1)M_{n-1} + x_n$



$\Rightarrow M_n = (1 - \frac{1}{n}) M_{n-1} + \frac{1}{n} x_n$

Online



$S_n/n = (1 - k_n) S_{n-1}$   
+  $k_n x_n$

# Reasons for Popularity of Kalman Filter

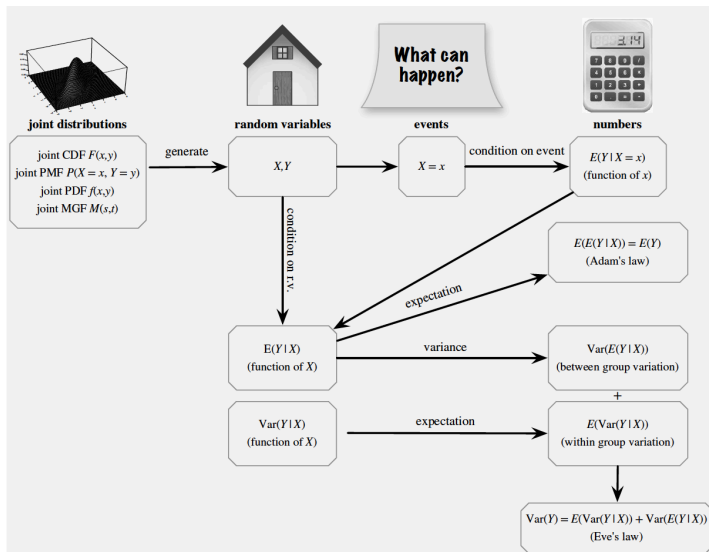
- Good results in practice due to optimality and structure: LLSE estimation in general, MMSE estimation under the setting of Gaussian noise.
- Convenient form for online real time processing: recursive equations.
- Easy to formulate and implement given a basic understanding.

# Why Use The Word “Filter”

- The process of finding the “best estimate” from noisy data amounts to “filtering out” the noise.
- Estimation (statistical perspective) vs. Filtering (signal processing perspective)
- A Kalman filter not only cleans up the data measurements
- A Kalman filter also projects these measurements onto the state estimate



# Summary 1



# References

- Chapter 9 of **BH**
- Chapters 4 & 6 of **BT**