

# Lecture 3: Conditional Expectation

Ziyu Shao

School of Information Science and Technology  
ShanghaiTech University

March 16,18,23,25, 2020

# Outline

- 1 Conditional Expectation Given An Event
- 2 Conditional Expectation Given An R.V.
- 3 Properties of Conditional Expectation
- 4 Application I: Prediction
- 5 Application II: Estimation
- 6 Application III: Branching Process
- 7 Application IV: Poisson Process
- 8 References

# Motivation

- 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:  $E(Y|X)$

# Outline

- 1 Conditional Expectation Given An Event
- 2 Conditional Expectation Given An R.V.
- 3 Properties of Conditional Expectation
- 4 Application I: Prediction
- 5 Application II: Estimation
- 6 Application III: Branching Process
- 7 Application IV: Poisson Process
- 8 References

# 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

$$E(Y|A) = \sum_y y P(Y = y|A),$$

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

$$E(Y|A) = \int_{-\infty}^{\infty} y f(y|A) dy,$$

where the conditional PDF  $f(y|A)$  is defined as the derivative of the conditional CDF  $F(y|A) = P(Y \leq y|A)$ , and can also be computed by a hybrid version of Bayes' rule:

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

# Direct Computing of $E(Y|A)$

①

$Y$  is continuous r.v.

$$\underbrace{P(Y \in (y-\varepsilon, y+\varepsilon) | A)} = \frac{P(A | Y \in (y-\varepsilon, y+\varepsilon))}{P(A)} P(Y \in (y-\varepsilon, y+\varepsilon))$$

$$\approx \cancel{\frac{1}{2\varepsilon} f(y|A)} = \frac{\cancel{P(A | Y \in (y-\varepsilon, y+\varepsilon))} \cdot \cancel{\frac{1}{2\varepsilon} f(y)}}{P(A)}$$

L.H.  
 $\varepsilon \rightarrow 0$ .

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

②  $E(Y|A) = \int_{-\infty}^{\infty} y f(y|A) dy$ ,  $y \in (-\infty, \infty)$

# Intuition for $E(Y|A)$

## Principle

$E(Y|A)$  is approximately the average of  $Y$  in a large number of simulation runs in which  $A$  occurred.  
*Monte Carlo..*

# Life Expectancy

①  $T$ : life span.

$$E(T) = 70 = E(T | T \geq 20)$$

$$\underline{E(T | T \geq 20)}$$

- Suppose you are 20 years old now.
- Average life time of your country is 70 years old.
- So you have 50 years to live on average. True or not?

② In general,  $E(T) \neq E(T | T \geq 20)$

③  $T \sim \text{exp}(\lambda)$ . memoryless.

# Law of Total Expectation $\text{LTE} \rightarrow \text{LOTE}$ .

$Y = I_B$  (indicator of event  $B$  occurring).

$$\begin{aligned} P(B) &= E(I_B) = E(Y) \stackrel{\text{LTE}}{=} \sum_{i=1}^n E(Y|A_i) P(A_i) = \sum_{i=1}^n E(I_B|A_i) P(A_i) \\ &= \sum_{i=1}^n P(B|A_i) P(A_i) \end{aligned}$$

## Theorem

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

$$E(Y) = \underbrace{\sum_{i=1}^n E(Y|A_i) P(A_i)}_{\text{LTE}}$$

# Outline

- 1 Conditional Expectation Given An Event
- 2 Conditional Expectation Given An R.V.
- 3 Properties of Conditional Expectation
- 4 Application I: Prediction
- 5 Application II: Estimation
- 6 Application III: Branching Process
- 7 Application IV: Poisson Process
- 8 References

# Conditional Expectation Given An R.V.

## 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 computer  $E(E(Y|X))$  and  $\text{Var}(E(Y|X))$ .

## Example

①  $M = \max(X, Y)$ ,  $L = \min(X, Y)$

By memoryless property,  $M-L$  is independent of  $L$

$$M-L \sim \text{expo}(\lambda)$$

②  $E(M|L)$

For  $X, Y \stackrel{\text{i.i.d.}}{\sim} \text{Expo}(\lambda)$ , find  $E(\max(X, Y) | \min(X, Y))$ .

$$\begin{aligned} E(M|L=c) &= E(M-L+L|L=c) \\ &= E(M-L|L=c) + E(L|L=c) \\ &= E(M-L) + c \\ &= \frac{1}{\lambda} + c = g(c) \end{aligned}$$

③  $E(M|L)$  =  $\frac{1}{\lambda} + L$

# Solution

# Outline

- 1 Conditional Expectation Given An Event
- 2 Conditional Expectation Given An R.V.
- 3 Properties of Conditional Expectation
- 4 Application I: Prediction
- 5 Application II: Estimation
- 6 Application III: Branching Process
- 7 Application IV: Poisson Process
- 8 References

# Dropping What's Independent

## Theorem

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

# Taking Out What's Known

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

$$= \underbrace{h(x)}_{\text{constant}} E[Y|X=x]$$

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

Theorem

For any function  $h$ ,

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

# Linearity

$$E(Y|X_1+X_2) \neq E(Y|X_1) + E(Y|X_2)$$

## Theorem

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

# Example

By symmetry property,  $E[X_1|S_n] = E[X_2|S_n] = \dots = E[X_n|S_n]$

By Linearity,  $E[X_1|S_n] + \dots + E[X_n|S_n]$

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

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

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

# Adam's Law

W.L.O.G. Focus on the scenario that  $X$  and  $Y$  are both discrete.

$$1^\circ. \quad g(X) = E(Y|X)$$

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

## Theorem

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

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

$$\begin{aligned} 2^\circ. \quad & E[g(X)] = \sum_x g(x) P(X=x) \\ &= \sum_X \left( \sum_y y P(Y=y|X=x) \right) \cdot P(X=x) \\ &= \underbrace{\sum_y y \sum_X P(Y=y|X=x) \cdot P(X=x)}_{\text{ }} = \sum_y y \cdot \underbrace{\sum_X P(Y=y, X=x)}_{\text{ }} \\ &= \sum_y y \cdot P(Y=y) = E(Y) \end{aligned}$$

# Proof

# Adam's Law and LOTE

$$\text{LOTE: } E(Y) = \sum_x E[Y|X=x] P(X=x)$$

Adam's Law

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

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

Since

$$E[E(Y|X)] = E[g(X)] = \sum_x g(x) P(X=x)$$

$$= \sum_x E[Y|X=x] P(X=x)$$

$$= E(Y)$$

# Adam's Law with Extra Conditioning

$$\hat{E}(\cdot) = E(\cdot | z) \quad \text{conditional expectation is expectation.}$$

Adam's Law.

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

$\hat{E}(Y|X, Z)$        $E(Y|Z)$

Theorem

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

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

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

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

important!  
Bandit / RL.  
Bellman equation.

# Conditional Variance

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

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

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

## Definition

The *conditional variance of Y given X* is

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

This is equivalent to

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

# Eve's law

## Theorem

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

$$\text{Var}(Y) = \underbrace{E(\text{Var}(Y|X))}_{\text{EV}} + \underbrace{\text{Var}(E(Y|X))}_{\text{VE}}.$$

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** ①  $g(x) = E(Y|x)$ . By Adam's Law,  $E[g(x)] = E(Y)$

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

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

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

## Example: Random Sum

① Conditioning on  $N$ ,  $E(\underline{x}|N=n) = E(\sum_{j=1}^N x_j | N=n) = E(\sum_{j=1}^n x_j | N=n)$   
 $= E(\sum_{j=1}^n x_j) = \sum_{j=1}^n E(x_j) = n \cdot \mu.$

$$\Rightarrow E(x|N) = N \cdot \mu \quad \Rightarrow E(x) = E[E(x|N)] = E(N \cdot \mu) = \mu \cdot E(N)$$

A store receives  $\underline{N}$  customers in a day, where  $\underline{N}$  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  $\underline{N}$  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  $\text{Var}(N)$ .

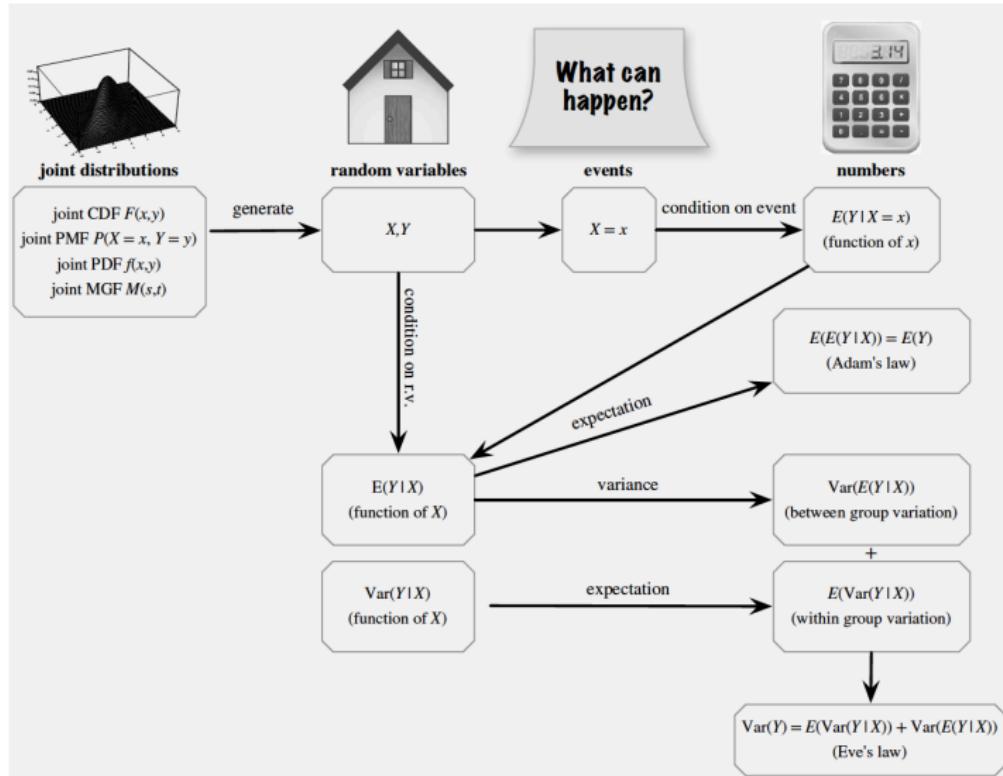
②  $\text{Var}(x|N=n) = \text{Var}(\sum_{j=1}^n x_j | N=n) = \text{Var}(\sum_{j=1}^n x_j | N=n)$   
 $= \text{Var}(\sum_{j=1}^n x_j) = \sum_{j=1}^n \text{Var}(x_j) = n \sigma^2$   
 $\Rightarrow \text{Var}(x|N) = N \cdot \sigma^2$

# Solution

③ By Eve's Law,  $\text{Var}(X) = E[\text{Var}(X|N)] + \text{Var}[E(X|N)]$

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

# Summary



# Outline

- 1 Conditional Expectation Given An Event
- 2 Conditional Expectation Given An R.V.
- 3 Properties of Conditional Expectation
- 4 Application I: Prediction
- 5 Application II: Estimation
- 6 Application III: Branching Process
- 7 Application IV: Poisson Process
- 8 References

# Prediction Perspective

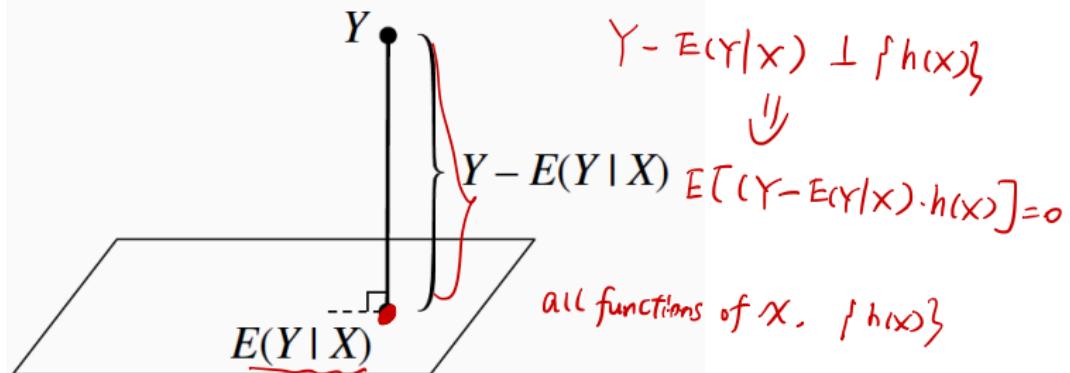
- Predict the future observations or unknown parameters based on data
- $E(Y|X)$  is our best predictor of  $Y$  based on data  $X$ .
- Best means it is the function of  $X$  with the lowest mean squared error (expected squared difference between  $Y$  and prediction of  $Y$ ).

# Geometric Perspective

1°. Projection of  $Y$  onto the plane  
of all functions of  $X$ .

$$a \perp b$$
$$E(a \cdot b) = 0$$

2°.  $Y - E(Y|X)$  is orthogonal to  
the plane.



# Projection Interpretation

## Theorem

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

$$\underline{E((Y - E(Y|X))h(X))} = 0.$$

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

Proof 1<sup>o</sup>.  $Z = Y - E(Y|X)$ ,  $E(Z) = E[Y - E(Y|X)] = E(Y) - E[E(Y|X)]$   
 $= E(Y) - E(Y) = 0$

$$\begin{aligned}\text{Cov}(Y - E(Y|X), h(X)) &= \text{Cov}(Z, h(X)) \\ &= E[(Z - EZ)(h(X) - E[h(X)])] \\ &= E[Z(h(X) - E[h(X)])] = E[Zh(X)] - E[h(X)] \cdot E(Z) \\ &= E[Zh(X)]\end{aligned}$$

thus  $\text{Cov}(Y - E(Y|X), h(X)) = 0 \Leftrightarrow E[Zh(X)] = 0$

$$\begin{aligned}2^o. \quad E[Zh(X)] &= E[(Y - E(Y|X))h(X)] \\ &= E(Yh(X)) - E[h(X)E(Y|X)] \\ &= E(Yh(X)) - E[E[Yh(X)|X]] \quad \text{taking back what's known} \\ &= E(Yh(X)) - E(Yh(X)) \quad \text{Adam's Law} \\ &= 0\end{aligned}$$

# Proof

# Prediction Perspective

$$\min_{\hat{Y}} E[(Y - \hat{Y})^2]$$

$$\begin{aligned}\hat{Y}^* &= \underset{\hat{Y}}{\operatorname{arg\,min}} E[(Y - \hat{Y})^2] \\ &= E(Y|X)\end{aligned}$$

MSE

Minimum Mean Squared Error.

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

Proof 1<sup>o</sup>.  $\hat{Y}$ : predictor of  $Y$ , a function of  $X$ .<sup>Info/data</sup>

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

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

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

2<sup>o</sup>. Let  $h(x) = E(Y|X) - g(x)$  is a function of  $X$ .  
By projection theorem, we know

$$E[(Y - E(Y|X)) \cdot h(x)] = 0$$

## Proof

$$\begin{aligned} 3^o. \quad & E[(Y - \hat{Y})^2] = E[(Y - g(x))^2] \\ &= \underbrace{E[(Y - E(Y|x))^2]}_{\geq 0} + \underbrace{E[(E(Y|x) - g(x))^2]}_{\geq 0} \\ &\geq E[(Y - E(Y|x))^2] \end{aligned}$$

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

$$\text{MMSE} \quad E[(Y - \hat{Y}^*)^2] = E[(Y - E(Y|x))^2]$$

# Proof

# Perspective of Statistical Learning

In general,  $E(Y|X)$  is a nonlinear function of  $X$ , which is hard to compute. At least half contents of statistical learning courses tell you how to approximately compute  $E(Y|X)$ .

- Linear regression: using linear function to approximate
- Logistic regression: using log-linear function to approximate
- Polynomial regression: using polynomial function to approximate
- Regression splines: using spline function to approximate
- Deep learning for prediction: using neural networks to approximate

# Outline

- 1 Conditional Expectation Given An Event
- 2 Conditional Expectation Given An R.V.
- 3 Properties of Conditional Expectation
- 4 Application I: Prediction
- 5 Application II: Estimation
- 6 Application III: Branching Process
- 7 Application IV: Poisson Process
- 8 References

# Basic Estimation Problem

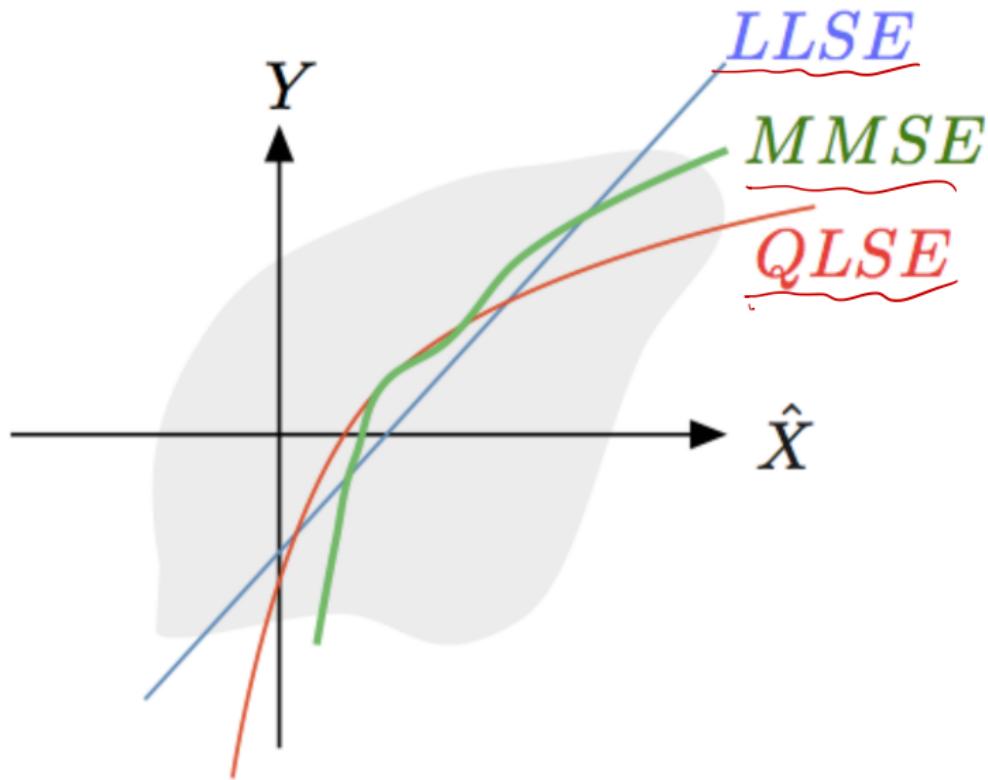
- Estimate  $X$  from the observed value  $Y$
- Choose the estimator (inference function)  $g(\cdot)$  to minimize the expected error  $E(c(X, g(Y)))$
- $c(X, \hat{X})$  is the cost of guessing  $\hat{X}$  when the actually value is  $X$ .

$$\hat{x} = g(x)$$

# Basic Estimation Problem

- When  $c(X, \hat{X}) = \|\mathbf{X} - \hat{\mathbf{X}}\|^2$ , the problem is called “the least square estimate (LSE)” problem.
- Further, if the function  $g(\cdot)$  is restricted to be linear, i.e., of the form  $a + bY$ , it is called “the Linear Least Square Estimate (LLSE)” estimate of  $X$  given  $Y$ .
- Further, if the function  $g(\cdot)$  is restricted to be quadratic, i.e., of the form  $a + bY^2$ , it is called “the Quadratic Least Square Estimate (QLSE)” estimate of  $X$  given  $Y$ .
- Further, if the function  $g(\cdot)$  can be arbitrary, it is called “the Minimum Mean Square Estimate (MMSE)” estimate of  $X$  given  $Y$ .

# Illustration



# Linear Least Square Estimate

## Theorem

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

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

Proof 1<sup>o</sup>. Cost function  $g(a,b) = E[(X-a-bY)^2]$

$$= a^2 - 2aE(X) + 2abE(Y) + b^2E(Y^2) - 2bE(XY) + E(X^2)$$

To minimize  $g(a,b)$ .  $\frac{\partial g(a,b)}{\partial a} = 0$ ,  $\frac{\partial g(a,b)}{\partial b} = 0$

$$2^o \quad \left. \frac{\partial g(a,b)}{\partial a} \right|_{\begin{matrix} a \\ b \end{matrix}}^* = 2a^* - 2E(X) + 2b^*E(Y) = 0$$

$$\left. \frac{\partial g(a,b)}{\partial b} \right|_{\begin{matrix} a \\ b \end{matrix}}^* = 2a^*E(Y) + 2b^*E(Y^2) - 2E(XY) = 0$$

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

$$3^o. L(X|Y) = a^* + b^*Y = E(X) - \frac{\text{Cov}(X,Y)}{\text{Var}(Y) \cdot E(Y)} + \frac{\text{Cov}(X,Y)}{\text{Var}(Y)} \cdot Y$$
$$= E(X) + \frac{\text{Cov}(X,Y)}{\text{Var}(Y)} \cdot (Y - E(Y))$$

# Proof

$$4^{\circ} \quad \min_{a,b} g(a,b), \quad \frac{\partial g(a,b)}{\partial a} = 0, \quad \frac{\partial g(a,b)}{\partial b} = 0$$

Hessian matrix.  $H = \begin{bmatrix} \frac{\partial^2 g(a,b)}{\partial a^2} & \frac{\partial^2 g(a,b)}{\partial a \partial b} \\ \frac{\partial^2 g(a,b)}{\partial b \partial a} & \frac{\partial^2 g(a,b)}{\partial b^2} \end{bmatrix} \succcurlyeq 0$

positive semidefinite.

$$\text{Let } \mathbf{z}^T = [c, d], \quad \mathbf{z}^T H \mathbf{z} \geq 0$$

$$H = \begin{bmatrix} 2 & 2E[Y] \\ 2E[Y] & 2E[Y^2] \end{bmatrix} = 2 \begin{bmatrix} 1 & E[Y] \\ E[Y] & E[Y^2] \end{bmatrix}$$

$$\Rightarrow \mathbf{z}^T H \mathbf{z} = 2[c, d] \begin{bmatrix} 1 & E[Y] \\ E[Y] & E[Y^2] \end{bmatrix} \begin{bmatrix} c \\ d \end{bmatrix}$$

$$= 2[(c + dE[Y])^2 + d^2 \text{var}(Y)]$$

$$\geq 0$$

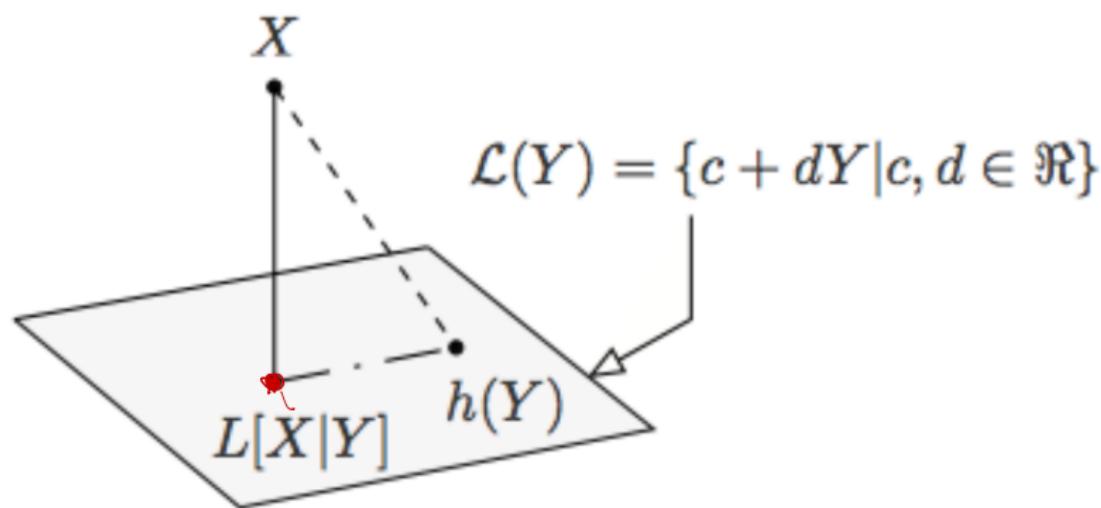
# Proof

# Projection Perspective

- Random variables  $X$  and  $Y$  are called “orthogonal” if  $E(XY) = 0$ , denoted by  $\underline{X \perp Y}$ .
- $E(\underline{X - E(X)}) = 0$ :  $\underline{X - E(X)} \perp c$  for any constant  $c$ .
- $E[(\underline{X - L(X|Y)})g(Y)] = 0$ :  $\underline{X - L(X|Y)} \perp \underline{g(Y)}$  for any linear functions of  $Y$ .

# Projection Perspective

$L[X|Y]$  is the projection of  $X$  onto the set  $\mathcal{L}(Y)$  of linear functions of  $Y$ :  $X - L[X|Y]$  is orthogonal to every linear function of  $Y$ .



# Linear Regression

Suppose we observe  $K$  i.i.d. samples  $(X_1, Y_1), \dots, (X_K, Y_K)$  of  $X, Y$ .  
Find a function  $g(Y) = a + bY$  such that  $E[(X - a - bY)^2]$  is minimized.

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

Monte Carlo

Sample average  $\rightarrow$  Expectation

# Linear Regression

$$L^k(x|Y) = E_k(x) + \frac{\text{Cov}_k(x, Y)}{\text{Var}_k(Y)} (Y - E_k(Y))$$

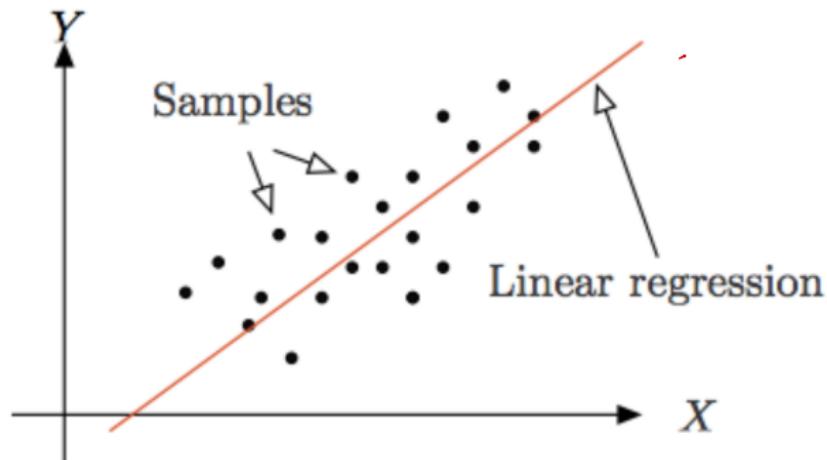
$$1^{\circ}. E_k(x) = \frac{1}{k} \sum_{j=1}^k x_j$$

$$2^{\circ}. E_k(y) = \frac{1}{k} \sum_{j=1}^k y_j$$

$$\begin{aligned} 3^{\circ}. \text{Cov}_k(x, Y) &= E_k(xy) - \underbrace{E_k(x) \cdot E_k(Y)} \\ &= \frac{1}{k} \sum_{j=1}^k x_j y_j - E_k(x) E_k(Y) \end{aligned}$$

$$\begin{aligned} 4^{\circ}. \text{Var}_k(Y) &= E_k(Y^2) - (E_k(Y))^2 \\ &= \frac{1}{k} \sum_{j=1}^k y_j^2 - (E_k(Y))^2 \end{aligned}$$

# Linear Regression



# Linear Regression

As  $k \rightarrow \infty$ , by SLLN.

$$E_k(x) \rightarrow E(x)$$

$$E_k(Y) \rightarrow E(Y)$$

$$\text{Cov}_k(X, Y) \rightarrow \text{Cov}(X, Y)$$

$$\text{Var}_k(Y) \rightarrow \text{Var}(Y) \quad \text{w.p.1}$$

$$\Rightarrow L^k(x|Y) \rightarrow L(x|Y)$$

## Theorem

As the number of samples increases, the linear regression approaches the LLSE.

# New Perspectives of Covariance

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

- 1<sup>o</sup>. if  $\text{Cov}(X,Y) = 0$ ,  $E[X|Y]$  does not depend on  $Y$
- 2<sup>o</sup>. if  $\text{Cov}(X,Y) > 0$   $E[X|Y] \uparrow$  with  $Y$  ( $Y - E[Y]$ )  
 $\downarrow$   $Y$  ( $Y - E[Y]$ )

$\text{Cov}(X,Y)$  measures a form of dependency in terms of LR.

# Minimum Mean Squared Error Estimator

## Theorem

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

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

# Orthogonality Property of MMSE

## Theorem

(a) For any function  $\phi(\cdot)$ , one has

$$\underbrace{E[(X - E[X|Y])\phi(Y)]}_{} = 0$$

(b) Moreover, if the function  $g(Y)$  is such that

$$\underbrace{E[(X - g(Y))\phi(Y)]}_{} = 0, \forall \phi(\cdot).$$

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

Proof (a).  $E[(x - E[x|Y])\phi(Y)]$

$$= E[x\phi(Y) - \phi(Y)E[x|Y]]$$

$$= E[x\phi(Y)] - \underbrace{E[\phi(Y)E[x|Y]]}$$

$$= E[x\phi(Y)] - \underbrace{E[E[x\phi(Y)|Y]]}_{\text{Adam's Law}}$$

$$= E[x\phi(Y)] - E[x\phi(Y)]$$

$$= 0$$

$$\therefore \phi(\cdot)$$

Proof

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

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

$$= \underbrace{E[(g(Y) - E(X|Y))(g(Y) - X)]}_{\Phi(Y) = g(Y) - E(X|Y)} + \underbrace{E[(g(Y) - E(X|Y)).(X - E(X|Y))]}_{X - E(X|Y) \perp \Phi(Y)}$$

a function of  $Y$ .

$g(Y) - X \perp \Phi(Y)$  by (b).

$X - E(X|Y) \perp \Phi(Y)$ .

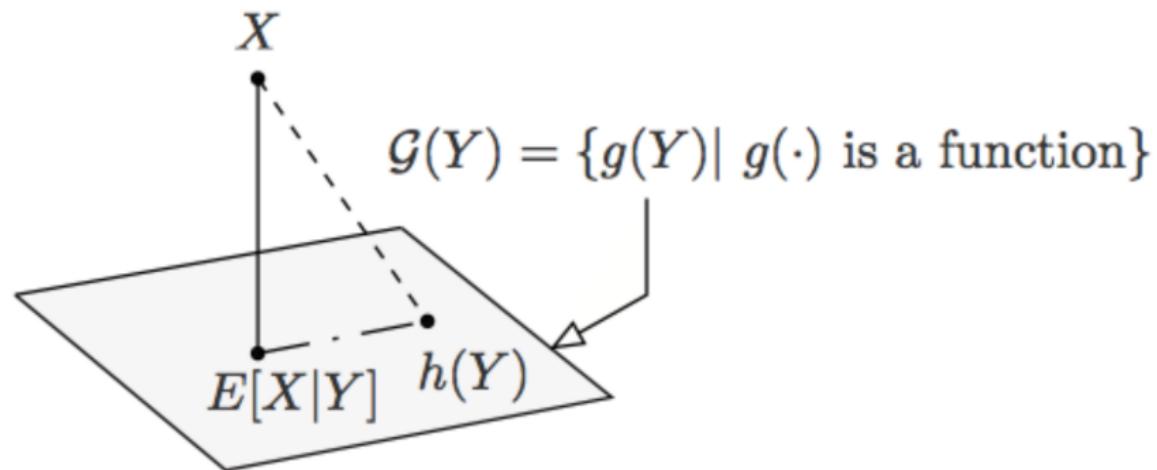
by (a).

$$= 0$$

$\nabla g(\cdot)$

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

# Projection Perspective of MMSE



# MMSE and LLSE

- In general, the MMSE estimator is different from the LLSE:  
 $E[X|Y] \neq L[X|Y]$ .
- The ordering of mean squared errors is

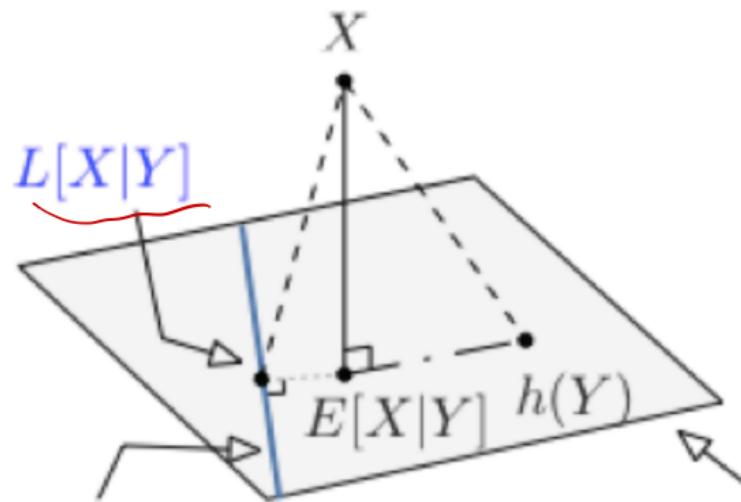
$$\underbrace{E[(X - E[X|Y])^2]}_{\text{MMSE.}} \leq \underbrace{E[(X - L[X|Y])^2]}_{\text{LLSE.}} \leq \underbrace{\text{Var}(X)}$$

$$E[(X - E(X))^2]$$

$E(X)$  as estimator.

a special case  
of linear  
function of  $X$ .

# MMSE and LLSE



$$\mathcal{L}(Y) = \{a + bY | a, b \in \mathbb{R}\}$$
$$\mathcal{G}(Y) = \{g(Y) | g(\cdot) \text{ is a function}\}$$

# MMSE for Jointly Gaussian Random Variables

## Theorem

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

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

## Example I

Let  $X, Y$  be jointly continuous random variables with the pdf

$$\underline{f_{X,Y}(x,y)} = \begin{cases} x + y & \text{if } 0 \leq x, y \leq 1 \\ 0 & \text{otherwise} \end{cases}$$

Find the MMSE  $E(X|Y)$  and LLSE  $L(X|Y)$ .

## Solution

1<sup>0</sup>. Find MMSE  $E[X|Y]$

$$f_Y(y) = \int_{-\infty}^{\infty} f_{X,Y}(x,y) dx = \begin{cases} \frac{1}{2} + y & 0 \leq y \leq 1 \\ 0 & \text{else} \end{cases}$$

$$f_{X|Y}(x|y) = \frac{f_{X,Y}(x,y)}{f_Y(y)} = \begin{cases} \frac{x+y}{\frac{1}{2}+y} & 0 \leq x \leq 1 \\ 0 & \text{else} \end{cases}$$

Thus for  $0 \leq y \leq 1$ ,  $E[X|Y=y] = \int_0^1 x f_{X|Y}(x|y) dx$

$$= \int_0^1 \frac{x^2 + xy}{\frac{1}{2} + y} dx = \frac{2 + 3y}{3 + 6y}$$

$$\Rightarrow E[X|Y] = \frac{2 + 3Y}{3 + 6Y}$$

# Solution

$$2^{\circ}. \text{ Find } \text{LSE } L[x|Y] = E[X] + \frac{\text{Cov}(X,Y)}{\text{Var}(Y)} (Y - E(Y))$$

$$E(X) = E(Y) = \frac{7}{12}$$

$$\text{Var}(Y) = \frac{11}{144} \quad , \quad \text{Cov}(X,Y) = -\frac{1}{144}$$

$$\Rightarrow L[x|Y] = \frac{7}{12} - \frac{1}{11} (Y - \frac{7}{12})$$

## Example II: Revisit Biased Coin Problem

We wish to estimate the probability of landing heads, denoted by  $\theta$ , of a biased coin. We model  $\theta$  as the value of a random variable  $\Theta$  with a known prior PDF  $f_\Theta \sim \text{Unif}(0, 1)$ . We consider  $n$  independent tosses and let  $X$  be the number of heads observed. Find the MMSE  $E(\Theta|X)$  and LLSE  $L(\Theta|X)$ .

# Solution

MMSE

1<sup>o</sup>.  $\Theta \sim \text{unif}(0,1) = \text{Beta}(1,1)$ , prior pdf

# of heads  $X | \Theta = \theta \sim \text{Bin}(n, \theta)$

By Beta-binomial Conjugacy,

$$\Theta | X=k \sim \text{Beta}(k+1, n-k+1)$$

$$\Rightarrow E[\Theta | X=k] = \frac{k+1}{n+2}$$

$$\Rightarrow E[\Theta | X] = \frac{n+1}{n+2}$$

MMSE

Univar of  $X$ .

$$E[\Theta | X] = L[\Theta | X]$$

# Solution

$$2^{\circ}. \text{LT}(\theta|x) = E[\theta] + \frac{\text{cov}(\theta, x)}{\text{var}(x)}(x - E[x])$$

$$\theta \sim \text{Unif}(0,1) \Rightarrow E[\theta] = \frac{1}{2}, \text{Var}(\theta) = \frac{1}{12}, E[\theta^2] = \frac{1}{3}.$$
$$x|\theta = \theta \sim \text{Bin}(n, \theta)$$

$$E[x] = E[E(x|\theta)] = E[n\theta] = n \cdot E[\theta] = \frac{n}{2}$$

$$\text{Var}(x) \stackrel{\text{Events Law}}{=} E[\text{Var}(x|\theta)] + \text{Var}(E(x|\theta))$$

$$= E[n\theta(1-\theta)] + \text{Var}(n\theta)$$

$$= n[E(\theta) - E(\theta^2)] + n^2 \text{Var}(\theta)$$

$$= n\left(\frac{1}{2} - \frac{1}{3}\right) + n^2 \cdot \frac{1}{12}$$

$$= \frac{n}{12}(n+2)$$

# Solution

$$\begin{aligned}3^{\circ}. \quad \text{Cov}(x, \theta) &= E(T\theta x) - E(\theta) \cdot E(x) \\&= E[E(T\theta x | \theta)] - E(\theta) \cdot E(x) \\&= E[T \theta E[x | \theta]] - E(\theta) \cdot E(x) \\&= E[\theta^2 \cdot n] - \frac{1}{2} \cdot \frac{n}{2} \\&= n E(\theta^2) - \frac{n}{4} \\&= n \cdot \frac{1}{3} - \frac{n}{4} = \frac{n}{12}\end{aligned}$$

$$\begin{aligned}4^{\circ} \quad L[\theta | x] &= E(\theta) + \frac{\text{Cov}(\theta, x)}{\text{var}(x)} [x - E(x)] \\&= \frac{1}{2} + \frac{\frac{n}{12}}{\frac{n}{12}(n+2)} [x - \frac{n}{2}] = \frac{x+1}{n+2} \\&= E(\theta | x)\end{aligned}$$

# Application Case: Kalman Filter

- 1960: Rudolph Emil Kalman (1930-2016) introduced what is known as Kalman filter.
- Milestones in Statistics & Signal Processing



# Reasons for Popularity of Kalman Filter

- Good results in practice due to optimality and structure: LLSE estimation in general and 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 clean up the data measurements
- A Kalman filter also projects these measurements onto the state estimate

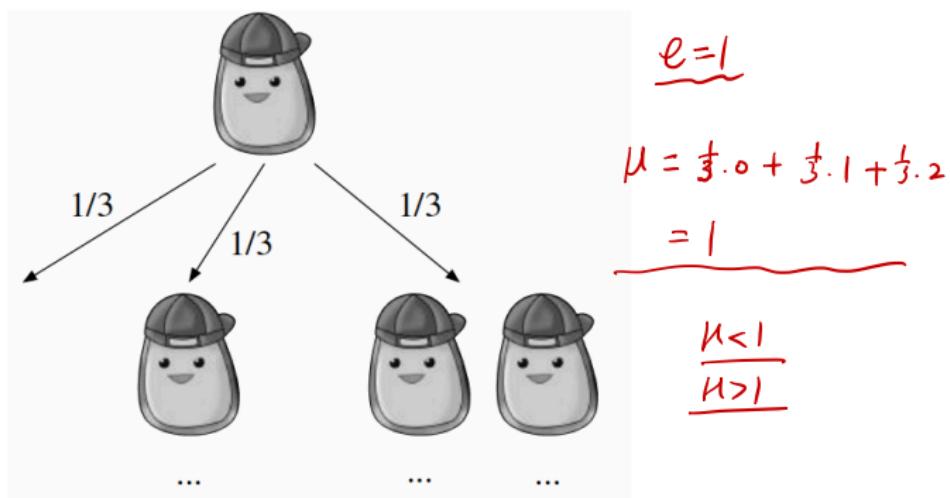
# Outline

- 1 Conditional Expectation Given An Event
- 2 Conditional Expectation Given An R.V.
- 3 Properties of Conditional Expectation
- 4 Application I: Prediction
- 5 Application II: Estimation
- 6 Application III: Branching Process
- 7 Application IV: Poisson Process
- 8 References

# Revisit the Story of Bobo

A single amoeba, Bobo, lives in a pond. After one minute Bobo will either die, split into two amoebas, or stay the same, with equal probability, and in subsequent minutes all living amoebas will behave the same way, independently. What is the probability that the amoeba population will eventually die out?

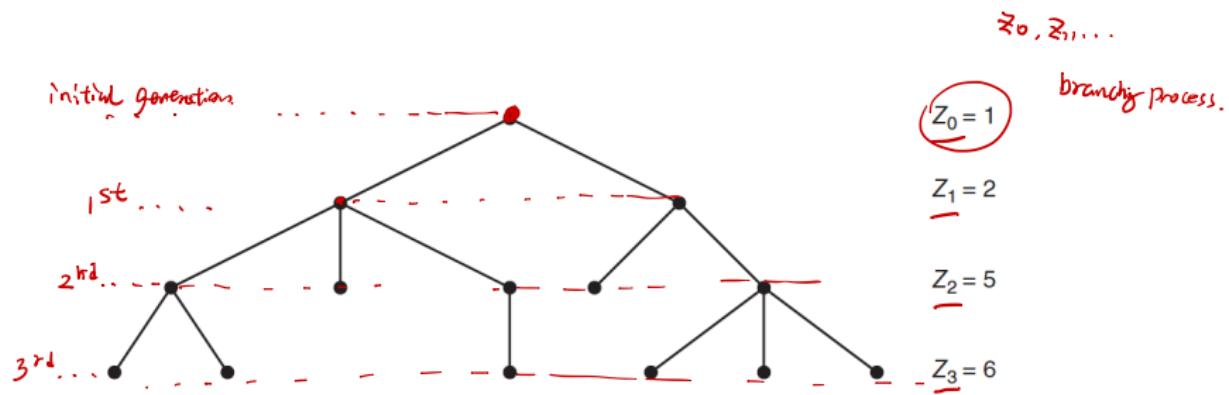
First Step Analysis.



# Branching Process

- Very useful model for the growth of populations
- Galton-Watson process: Sir Francois Galton & Reverend Henry William Watson made important contributions.
- Various Applications
  - ▶ Extinction of family surnames (original motivation)
  - ▶ The spread of infectious diseases and epidemics (biology and epidemiology)
  - ▶ The spread of computer software viruses
  - ▶ Rumor(information) spread in social networks

# Illustration of A Branching Process



offspring distribution  $a = (a_0, a_1, a_2, \dots)$

$a_0 = 1 \Rightarrow Z_n = 0, n \geq 1.$

Def:  $a_n$  prob. of an individual gives birth to  $k$  children.

Usually  $0 < a_0 < 1$ ,  $0 < a_0 + a_1 < 1$ ,  $\geq 2$  children w.p.  $> 0$ .

# Setting

- Offspring distribution  $(a) = (a_0, a_1, a_2, \dots)$ : distribution of each individual independently produces a random number of children
- An individual gives birth to  $k$  children with probability  $a_k$ ,  $k \geq 0$ , independently of other individuals.
- $Z_n$ : the size (number of individuals) of the  $n^{\text{th}}$  generation,  $n \geq 0$ .
- $Z_0, Z_1, \dots$  is a branching process.

# Mean Generation Size

①  $X_i$ : the  $i^{th}$  person's children number. In the  $(n-1)^{th}$  generation  
 $X_1, \dots, X_{n-1}, a$ .

$$E(X_i) = \mu, \quad X_i \text{ and } Z_{n-1} \text{ are independent}$$

## Theorem

Let  $\mu = \sum_{k=0}^{\infty} k \cdot a_k$ . Suppose  $Z_0 = 1$ , then the mean of  $n^{th}$  generation size is

$$E(Z_n) = \mu^n, n \geq 0.$$

$$\begin{aligned} ② \quad Z_n &= \sum_{i=1}^{Z_{n-1}} X_i \quad \Rightarrow E[Z_n | Z_{n-1}=k] = E\left[\sum_{i=1}^{Z_{n-1}} X_i | Z_{n-1}=k\right] \\ &= E\left[\sum_{i=1}^k X_i | Z_{n-1}=k\right] = E\left[\sum_{i=1}^k X_i\right] = k E[X_i] = k \cdot \mu \quad \Rightarrow E[Z_n | Z_{n-1}] = k Z_{n-1} \end{aligned}$$

$$\Rightarrow E[Z_n] = E[E[Z_n | Z_{n-1}]] = E[Z_{n-1} \mu] = \mu \cdot E[Z_{n-1}] \Rightarrow E(Z_n) = \mu^n, n \geq 0$$

# Proof

# Three Cases

$$E(Z_n) = \mu^n$$

$$\lim_{n \rightarrow \infty} E(Z_n) = \lim_{n \rightarrow \infty} \mu^n = \begin{cases} 0 & \text{if } \mu < 1 \\ 1 & \text{if } \mu = 1 \\ \infty & \text{if } \mu > 1 \end{cases}$$

Subcritical  
Critical  
Supercritical.

# Three Cases

TABLE 4.1 Simulations of a Branching Process for Three Choices of  $\mu$

$\mu$	$Z_0$	$Z_1$	$Z_2$	$Z_3$	$Z_4$	$Z_5$	$Z_6$	$Z_7$	$Z_8$	$Z_9$	$Z_{10}$
0.75	1	2	3	3	2	1	1	0	0	0	0
0.75	1	0	0	0	0	0	0	0	0	0	0
0.75	1	2	0	0	0	0	0	0	0	0	0
0.75	1	0	0	0	0	0	0	0	0	0	0
0.75	1	3	3	1	3	1	0	0	0	0	0
1	1	1	1	0	0	0	0	0	0	0	0
1	1	2	0	0	0	0	0	0	0	0	0
1	3	6	6	5	6	7	8	8	8	6	5
1	1	3	4	1	2	1	0	0	0	0	0
1	1	2	1	1	2	1	2	1	0	0	0
1.5	1	2	3	10	22	41	93	173	375	763	1,597
1.5	1	1	1	1	2	4	7	9	11	19	29
1.5	1	4	5	18	34	68	127	246	521	1,011	2,065
1.5	1	1	2	0	0	0	0	0	0	0	0
1.5	1	2	5	3	2	6	9	17	18	13	19

# Extinction in the Subcritical Case

## Theorem

*With probability 1, a subcritical branching process eventually goes extinct.*

# Variance of Generation Size

## Theorem

Let  $\sigma^2$  denote the variance of the offspring distribution. Then we have

$$\text{Var}(Z_n) = \begin{cases} n\sigma^2 & \text{if } \mu = 1 \\ \sigma^2 \mu^{n-1} \left( \frac{\mu^n - 1}{\mu - 1} \right) & \text{if } \mu \neq 1 \end{cases}$$

Proof

$$Z_n = \sum_{i=1}^{n-1} X_i \quad , \quad \text{Var}(Z_n) = E[\text{Var}(Z_n | Z_{n-1})] + \text{Var}(E(Z_n | Z_{n-1})) \quad n \geq 1$$

$$1^{\circ}. \quad E[Z_n | Z_{n-1} = k] = E\left(\sum_{i=1}^k X_i\right) = k \cdot \mu \Rightarrow E[Z_n | Z_{n-1}] = \mu \cdot Z_{n-1}$$

$$2^{\circ}. \quad \text{Var}(Z_n | Z_{n-1} = k) = \text{Var}\left(\sum_{i=1}^k X_i\right) = k \sigma^2 \Rightarrow \text{Var}(Z_n | Z_{n-1}) = \sigma^2 Z_{n-1}$$

$$\Rightarrow \text{Var}(Z_n) = E(\text{Var}(Z_n | Z_{n-1})) + \text{Var}(E(Z_n | Z_{n-1}))$$

$$= E(\sigma^2 Z_{n-1}) + \text{Var}(\mu Z_{n-1})$$

$$= \sigma^2 E(Z_{n-1}) + \mu^2 \text{Var}(Z_{n-1}) = \sigma^2 \cdot \mu^{n-1} + \mu^2 \text{Var}(Z_{n-1}) \quad n \geq 1$$

(a)  $\mu = 1 \Rightarrow \text{Var}(Z_n) = \text{Var}(Z_{n-1}) + \sigma^2 \quad n \geq 1 \Rightarrow \text{Var}(Z_n) = n \sigma^2, \quad n \geq 1$

(b)  $\mu \neq 1 \Rightarrow \text{Var}(Z_1) = \sigma^2, \quad \text{Var}(Z_2) = \sigma^2 \mu (\mu + 1)$

$$\text{Var}(Z_3) = \sigma^2 \mu^2 (\mu + \mu + \mu^2) \dots \text{Var}(Z_n) = \sigma^2 \mu^{n-1} \left( \sum_{k=0}^{n-1} \mu^k \right)$$

$$= \sigma^2 \mu^{n-1} \cdot \frac{\mu^n - 1}{\mu - 1} \quad n \geq 1$$

# Proof

$$E(Z_n) = \mu^n$$

$$\text{Var}(Z_n) = \begin{cases} n\sigma^2 & \mu=1 \\ \sigma^2 \mu^{n-1} \frac{\mu^n - 1}{\mu - 1} & \mu \neq 1 \end{cases}$$

$\mu < 1$  Subcritical.  $E(Z_n) \rightarrow 0$ ,  $\text{Var}(Z_n) \rightarrow 0$

$\mu = 1$  Critical.  $E(Z_n) = 1$ ,  $\text{Var}(Z_n) = n\sigma^2 \rightarrow \infty$

$\mu > 1$  Supercritical.  $E(Z_n) > 1$ ,  $\text{Var}(Z_n) \rightarrow \infty$

z

# Probability Generation Function (PGF)

## Theorem

*The PGF for  $Z_n$  is the  $n$ -fold composition of the PGF for offspring distribution.*

# Proof

1°.  $Z_n$ . PGF of  $Z_n$ :  $G_{n|S} = E[S^{Z_n}] = \sum_{k=0}^{\infty} S^k P(Z_n=k)$

On the other hand, let  $G_{1|S} = \sum_{k=0}^{\infty} S^k k a_k = E(S^X)$  be the PGF  
of the offspring distribution,

$$\Rightarrow G_{n|S} = E(S^{Z_n}) = E(S^{\sum_{j=1}^{Z_{n-1}} X_j}) = E\left[E\left(S^{\sum_{j=1}^{Z_{n-1}} X_j} | Z_{n-1}\right)\right]$$

independence of  $X_j$  and  $Z_{n-1}$

$$\begin{aligned} E(S^{\sum_{j=1}^{Z_{n-1}} X_j} | Z_{n-1}=z) &= E(S^{\sum_{j=1}^z X_j} | Z_{n-1}=z) = E(S^{\sum_{j=1}^z X_j}) \\ &= E\left(\prod_{j=1}^z S^{X_j}\right) = \prod_{j=1}^z E(S^{X_j}) = [G_{1|S}]^z \quad \text{by } z, (x_1, \dots, x_z) \sim i.i.d. \\ \Rightarrow E(S^{\sum_{j=1}^{Z_{n-1}} X_j} | Z_{n-1}) &= [G_{1|S}]^{Z_{n-1}} \end{aligned}$$

Take expectation.  $G_{n|S} = E[S^{Z_n}] = E[(G_{1|S})^{Z_{n-1}}] = G_{n-1}(G_{1|S})$

$$\Rightarrow G_{n|S} = (G_{n-1}(G_{1|S})) \quad \begin{cases} G_{0|S} = E[S^{Z_0}] = S, \quad G_{1|S} = G_0(G_{1|S}) = G_{1|S} \\ G_2(S) = G_1(G_{1|S}) = G[G_{1|S}] = G[G_{1|S}] \\ \dots \end{cases}$$

# Proof

$$G_n(S) = G[\dots G[G[S] \dots]] = G[G_{n+1}(S)]$$

*n-fold*

# Extinction Probability

$$\ell = \lim_{n \rightarrow \infty} e_n = \lim_{n \rightarrow \infty} p(Z_n=0)$$

$$e_n = p(Z_n=0) = G(e_{n-1}) = G(G(e_{n-1})) = G^2(e_{n-1})$$

$$\lim_{n \rightarrow \infty} e_n, \Rightarrow \ell = G(\ell) \text{ fixed point.}$$

## Theorem

Given a branching process, let  $G$  be the probability generating function of the offspring distribution. Then, the probability of eventual extinction is the smallest positive root of the equation  $s = G(s)$ . Further, If  $\mu \leq 1$ , that is, in the subcritical and critical cases, the extinction probability is equal to 1.

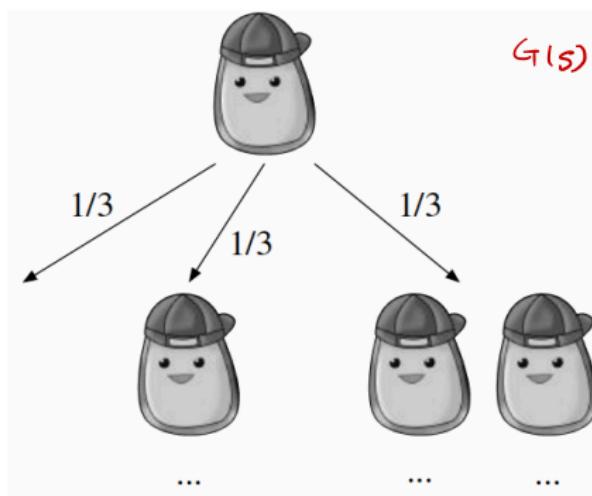
## Example: Story of Bobo

A single amoeba, Bobo, lives in a pond. After one minute Bobo will either die, split into two amoebas, or stay the same, with equal probability, and in subsequent minutes all living amoebas will behave the same way, independently. What is the probability that the amoeba population will eventually die out?

$$\alpha = (0, 1, 2)$$

$$\frac{1}{3} \quad \frac{1}{3} \quad \frac{1}{3}$$

$$G(s) = \sum_{k=0}^{\infty} s^k a_k = \frac{1}{3} + \frac{1}{3}s + \frac{1}{3}s^2$$



$$s = G(s)$$

$$\Rightarrow s^2 - 2s + 1 = 0$$

$$\Rightarrow s = 1$$

$$\Rightarrow \underline{s = 1}$$

# Example: First-Success Distribution of Offspring

$$1^{\circ}. \quad G(s) = \sum_{k=0}^{\infty} s^k a_k = \sum_{k=0}^{\infty} s^k (1-p)^k p = \frac{p}{1-s(1-p)} \quad (s(1-p)<1).$$

$$\mu = G'(1) = \frac{p}{1-p}. \quad \text{Supercritical} \Rightarrow \mu > 1 \Rightarrow p < \frac{1}{2}.$$

$$G(s) = a_0 + \sum_{k=1}^{\infty} s^k a_k$$

$$G'(s) = \sum_{k=1}^{\infty} k s^{k-1} a_k \Rightarrow G'(1) = \sum_{k=1}^{\infty} k a_k = \sum_{k=0}^{\infty} k a_k = \mu$$

A branching process has offspring distribution  $a_k = (1-p)^k p$ , for  $k = 0, 1, \dots$ . Find the extinction probability in the supercritical case.

2<sup>o</sup>. Solve the equation

$$S = G(s) = \frac{p}{1-s(1-p)}$$

$$\Rightarrow (1-p)s^2 - s + p = 0$$

$$\Rightarrow S_1 = 1, S_2 = \frac{p}{1-p} \quad \text{Since } \alpha < 1$$

$$\Rightarrow S_2 < S_1$$

$$\Rightarrow e = S_2 = \frac{p}{1-p}$$

# Simulation Results

$$P = \frac{1}{4}, \quad \mu = \frac{1+P}{P} = 3, \quad e = \frac{\mu}{\mu - P} = \frac{1}{2}.$$
$$\alpha_k = \left(\frac{3}{4}\right)^k \cdot \frac{1}{2}.$$

TABLE 4.2 Simulation of a Supercritical Branching Process, with  $\mu = 3$ . Four of the 12 runs go extinct by the 10th generation.

Z <sub>0</sub>	Z <sub>1</sub>	Z <sub>2</sub>	Z <sub>3</sub>	Z <sub>4</sub>	Z <sub>5</sub>	Z <sub>6</sub>	Z <sub>7</sub>	Z <sub>8</sub>	Z <sub>9</sub>	Z <sub>10</sub>
1	4	25	97	394	1160	3475	10685	31885	95757	287130
1	1	1	3	7	11	47	165	515	1525	4689
1	10	37	115	350	1124	3455	10073	29896	88863	267386
1	1	1	2	2	3	0	0	0	0	0
1	0	0	0	0	0	0	0	0	0	0
1	8	31	71	248	779	2282	6864	19895	59196	178171
1	0	0	0	0	0	0	0	0	0	0
1	4	7	13	34	106	380	1123	3385	10200	30090
1	16	49	163	447	1284	3794	11592	34626	104390	312704
1	1	1	5	16	51	155	559	1730	5378	15647
1	1	3	31	79	267	883	2637	8043	23970	71841
1	0	0	0	0	0	0	0	0	0	0

# Outline

- 1 Conditional Expectation Given An Event
- 2 Conditional Expectation Given An R.V.
- 3 Properties of Conditional Expectation
- 4 Application I: Prediction
- 5 Application II: Estimation
- 6 Application III: Branching Process
- 7 Application IV: Poisson Process
- 8 References

# Poisson Process

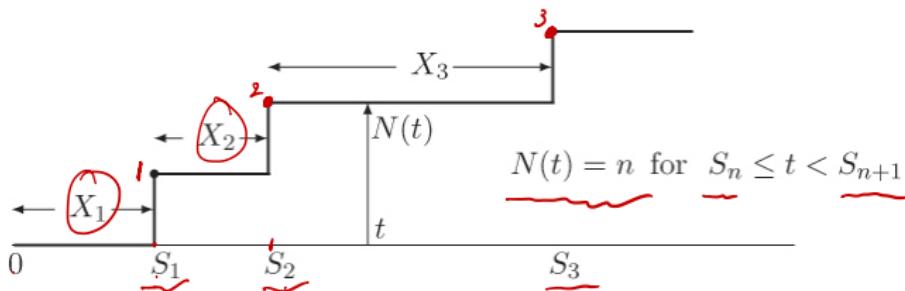


Figure : Illustration of a Poisson process and its arrival times  $\{S_1, S_2, \dots\}$ , its interarrival times  $\{X_1, X_2, \dots\}$ , and its counting process  $N(t) : t \geq 0$ .

# Various Perspectives

The Poisson process can be described in three ways:

- The interarrival times  $X_1, X_2, X_3, \dots$
- The arrival times  $S_0 = 0, S_1, S_2, S_3, \dots$
- The arrival counts  $N(t)$  in  $[0, t]$ ,  $t \geq 0$

## Important Relationships

$$S_{n-1} - S_{n-1} + S_{n-1} - S_{n-2} + \dots + S_1 - S_0.$$

$$S_n = \sum_{j=1}^n X_j, n \geq 1$$

$$X_n = S_n - S_{n-1}, n \geq 1$$

$$N(t) = \max\{n : S_n \leq t\}$$

$$\underline{N(t) \geq n} \Leftrightarrow \underline{S_n \leq t}$$

# Equivalent Definitions of Poisson Process



# Poisson Process - Definition 1

## Definition

A *Poisson process* with parameter  $\lambda$  is a counting process  $(N_t)_{t \geq 0}$  with the following properties:

- ①  $N_0 = 0$ .
- ② For all  $t > 0$ ,  $N_t$  has a Poisson distribution with parameter  $\lambda t$ .
- ③ (Stationary increments) For all  $s, t > 0$ ,  $N_{t+s} - N_s$  has the same distribution as  $N_t$ . That is,

$$P(N_{t+s} - N_s = k) = P(\underline{N_t} = k) = \frac{e^{-\lambda t} (\lambda t)^k}{k!}, \text{ for } k = 0, 1, \dots$$

- ④ (Independent increments) For  $0 \leq q < r \leq s < t$ ,  $\overset{\circ}{N}_t - \overset{\circ}{N}_s$  and  $\overset{\circ}{N}_r - \overset{\circ}{N}_q$  are independent random variables.

# Translated Poisson Process

Let  $(N_t)_{t \geq 0}$  be a Poisson process with parameter  $\lambda$ . For  $s > 0$ , let  $\tilde{N}_t = N_{t+s} - N_s$ , for  $t \geq 0$ . Then we have

- $(\tilde{N}_t)_{t \geq 0}$  is called “Translated Poisson Process”.
- $(\tilde{N}_t)_{t \geq 0}$  is a Poisson process with parameter  $\lambda$

# Poisson Process - Definition 2

## Definition

Let  $X_1, X_2, \dots$  be a sequence of *i.i.d.* exponential random variables with parameter  $\lambda$ . For  $t > 0$ , let

$$\underline{N_t = \max\{n : X_1 + \cdots + X_n \leq t\}},$$

with  $N_0 = 0$ . Then,  $(N_t)_{t \geq 0}$  defines a Poisson process with parameter  $\lambda$ .

# Arrival Times & Gamma Distribution

$$S_n = \underbrace{x_1 + \dots + x_n}_{X_1, \dots, X_n \text{ i.i.d.}} \sim \text{exp}(\lambda).$$

For  $n = 1, 2, \dots$ , let  $S_n$  be the time of the  $n$ -th arrival in a Poisson process with parameter  $\lambda$ . Then,  $S_n$  has a gamma distribution with parameters  $n$  and  $\lambda$ . The density function of  $S_n$  is

$$f_{S_n}(t) = \frac{\lambda^n t^{n-1} e^{-\lambda t}}{(n-1)!}, \text{ for } t > 0.$$

Mean and variance are

$$\underbrace{E(S_n)}_{\frac{n}{\lambda}} \text{ and } \underbrace{Var(S_n)}_{\frac{n}{\lambda^2}}$$

$$\begin{aligned} E(S_n) &= E(x_1) + \dots + E(x_n) \\ &= \frac{1}{\lambda} \cdot n = \frac{n}{\lambda}. \end{aligned}$$

# One New Perspective

CDF  $\xrightarrow{'} \text{PDF}$

Find  $f_{S_n}(t)$

$$i^{\circ} \cdot N(t) \geq n \Leftrightarrow S_n \leq t$$

$\Leftrightarrow$  CDF of  $S_n$

$$F_{S_n}(t) = P(S_n \leq t) = P(N(t) \geq n)$$

$$= \sum_{j=n}^{\infty} P(N(t)=j) = \sum_{j=n}^{\infty} \frac{e^{-\lambda t} (\lambda t)^j}{j!}$$

$$\text{PDF } f_{S_n}(t) = F'_{S_n}(t) = \left( \sum_{j=n}^{\infty} \frac{e^{-\lambda t} (\lambda t)^j}{j!} \right) /$$

= . . .

$$= \frac{\lambda e^{-\lambda t} (\lambda t)^{n-1}}{(n-1)!}$$

# Another New Perspective

$$\xrightarrow{s_{n+1} \quad t \quad s_n \quad t+dt}$$

$$f_{S_n}(t) dt = P(t < S_n < t+dt) = P(N(t)=n-1, \text{arrival in } [t, t+dt])$$

Independent increment

$$= P(N(t)=n-1) \cdot P\{\text{One arrival in } [t, t+dt]\}$$

stationary increment

$$= \frac{e^{-\lambda t} (\lambda t)^{n-1}}{(n-1)!} \cdot P\{N(dt)=1\}$$

$$= \frac{e^{-\lambda t} (\lambda t)^{n-1}}{(n-1)!} \cdot e^{-\lambda dt} \cdot \lambda dt$$

$$\Rightarrow f_{S_n}(t) = \lim_{dt \rightarrow 0} \frac{P(t < S_n < t+dt)}{dt} = \lim_{dt \rightarrow 0} \frac{e^{-\lambda t} (\lambda t)^{n-1}}{(n-1)!} \cdot \lambda e^{-\lambda dt}$$

$$= \frac{\lambda e^{-\lambda t} (\lambda t)^{n-1}}{(n-1)!}$$

# Arrival Times & Uniform Distribution

Let  $S_1, S_2, \dots$ , be the arrival times of a Poisson process with parameter  $\lambda$ . Conditional on  $N_t = n$ , the joint distribution of  $(S_1, \dots, S_n)$  is the distribution of the order statistics of  $n$  i.i.d. uniform random variables on  $[0, t]$ . That is, the joint density function of  $S_1, \dots, S_n$  is

$$f(s_1, \dots, s_n | N_t = n) = \frac{n!}{t^n}, \text{ for } 0 < s_1 < \dots < s_n < t.$$

Equivalently, let  $U_1, \dots, U_n$  be an i.i.d. sequence of random variables uniformly distributed on  $[0, t]$ . Then, conditional on  $N_t = n$ ,

$$(S_1, \dots, S_n) \text{ and } (U_{(1)}, \dots, U_{(n)})$$

have the same distribution.

## Example

1°.  $N(t)$ : # of people arriving by the time  $t$ .

$$N(t) \sim \text{Pois}(\lambda t)$$

2°. total waiting time of people who arrive before the band starts.

$$\sum_{k=1}^{N(t)} (t - S_k).$$

Concert-goers arrive at a show according to a Poisson process with parameter  $\lambda$ . The band starts playing at time  $t$ . The  $k$ -th person to arrive in  $[0, t]$  waits  $t - S_k$  time units for the start of the concert, where  $S_k$  is the  $k$ -th arrival time. Find the expected total waiting time of concert-goers who arrive before the band starts.

3°. we focus.  $E\left[\sum_{k=1}^{N(t)} (t - S_k)\right]$

$$\text{Solution}^{\text{d}}. E\left(\sum_{k=1}^{N(t)} (t-s_k)\right) = E\left[E\left[\sum_{k=1}^{N(t)} (t-s_k) \mid N(t)\right]\right]$$

$$\text{Consider } E\left[\sum_{k=1}^{N(t)} (t-s_k) \mid N(t)=n\right] = E\left[\sum_{k=1}^n (t-s_k) \mid N(t)=n\right]$$

Given  $N(t)=n$ ,  $(s_1, \dots, s_n) \sim (U_{(1)}, \dots, U_{(n)})$

$$\begin{aligned} \Rightarrow E\left(\sum_{k=1}^n (t-s_k) \mid N(t)=n\right) &= E\left(\sum_{k=1}^n (t-U_{(k)}) \mid N(t)=n\right) \\ &= E\left(t_n - \sum_{k=1}^n U_{(k)} \mid N(t)=n\right) \quad \left(\sum_{k=1}^n U_{(k)} = \sum_{k=1}^n U_k\right) \\ &= E\left(t_n - \sum_{k=1}^n U_k \mid N(t)=n\right) \quad \text{($U_1, \dots, U_n$ i.i.d. unif[0,t])} \\ &= E[t_n - \sum_{k=1}^n U_k] = t_n - \sum_{k=1}^n E(U_k) = t_n - n \cdot \frac{t}{2} = \frac{1}{2}tn \end{aligned}$$

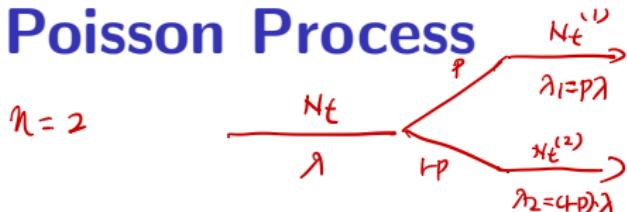
# Solution

We have  $E \left( \sum_{k=1}^{N(t)} (t-s_k) \mid N(t)=n \right) = \frac{1}{2} t \cdot n \quad \forall n$

$$\Rightarrow E \left( \sum_{k=1}^{N(t)} (t-s_k) \mid N(t) \right) = \frac{1}{2} t N(t).$$

$$\begin{aligned}\Rightarrow E \left( \sum_{k=1}^{N(t)} (t-s_k) \right) &= E \left[ E \left[ \sum_{k=1}^{N(t)} (t-s_k) \mid N(t) \right] \right] \\ &= E \left[ \frac{1}{2} t N(t) \right] \\ &= \frac{1}{2} t E[N(t)] \\ &= \frac{1}{2} t \lambda t \\ &= \frac{1}{2} \lambda t^2\end{aligned}$$

# Thinned Poisson Process

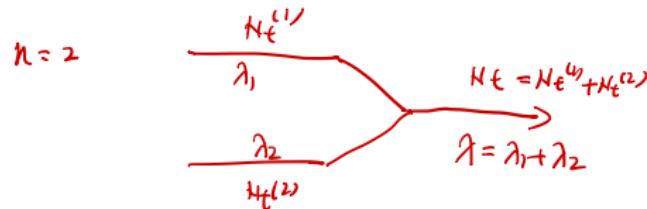


Let  $(N_t)_{t \geq 0}$  be a Poisson process with parameter  $\lambda$ . Assume that each arrival, independent of other arrivals, is marked as a type  $k$  event with probability  $p_k$ , for  $k = 1, \dots, n$ , where  $p_1 + \dots + p_n = 1$ . Let  $N_t^{(k)}$  be the number of type  $k$  events in  $[0, t]$ . Then,  $(N_t^{(k)})_{t \geq 0}$  is a Poisson process with parameter with  $\lambda p_k$ . Furthermore, the processes

$$(N_t^{(1)})_{t \geq 0}, \dots, (N_t^{(n)})_{t \geq 0}$$

are independent. Each process is called a *thinned Poisson process*.

# Superposition Process



Assume that  $(N_t^{(1)})_{t \geq 0}, \dots, (N_t^{(n)})_{t \geq 0}$  are  $n$  independent Poisson processes with respective parameters  $\lambda_1, \dots, \lambda_n$ . Let  $N_t = \underbrace{N_t^{(1)} + \dots + N_t^{(n)}}_{\text{for } t \geq 0}$ . Then,  $(N_t)_{t \geq 0}$  is a Poisson process with parameter  $\lambda = \lambda_1 + \dots + \lambda_n$ .

# Revisited Coupon Collector Problem

Embedding Trick. (Poisson process  $N(t)$ ,  $\lambda=1$ )

Discrete - Continuous.

- Previously we introduce the coupon collector problem:  $n$  toy types, collected one by one, sampling with replacement from the set of toy types each time, and all toy types were equally likely to be collected.
- Now suppose that at each stage, the  $j^{th}$  toy type is collected with probability  $p_j$ , where the  $p_j, j = 1, \dots, n$  are not necessarily equal.
- Let  $N$  be the number of toys needed until we have a full set.
- Find  $E(N)$ .

# Solution

(a). Suppose that the toy arrive according to a Poisson process with rate  $\lambda = 1$ .

Interarrival times between toys are iid.  $X_j \sim \text{expo}(1)$ ,  $j=1, 2, \dots, n$

Let  $Y_j$  = the waiting time until the first toy of type  $j$

[ By thinning the Poisson process  $\Rightarrow n$  independent Poisson process.

$[0, t]$  # of arrives type- $j$  toys.  $N_j(t) \sim \text{Pois}(\lambda p_j) = \text{Pois}(p_j t)$  ]

$$\Rightarrow Y_j \sim \text{expo}(p_j)$$

(b). The waiting time  $T$  until all toy types are collected as

$$T = \max(Y_1, \dots, Y_n)$$

$N$ : # of toys collected,  $T = X_1 + X_2 + \dots + X_N$  (Pois(1))  
Cumulative time of  $N$ th toys

$$\Rightarrow E(T) = E(X_1 + \dots + X_N) = E[E(X_1 + \dots + X_N | N)]$$

$$E[X_1 + \dots + X_N | N=n] = E[X_1 + \dots + X_n | N=n] = n E[X_1 | N=n] = n E[X_1] = n \cdot 1/n$$

$$\Rightarrow E[X_1 + \dots + X_N | N] = N$$

# Solution

$$\Rightarrow E(T) = E[E(X_1 + \dots + X_n | N)] = E(N)$$

Continuous  $\leftarrow$  bridge  $\xrightarrow{\text{poly}(\nu)} \text{discrete.}$

$$(c). E(N) = E(T) = \int_0^\infty P(T > t) dt$$

$$P(T > t) = 1 - P(T \leq t) = 1 - P(\max(Y_1, \dots, Y_n) \leq t)$$

$$= 1 - P(Y_1 \leq t, \dots, Y_n \leq t) = 1 - P(Y_1 \leq t) \dots P(Y_n \leq t) \quad Y_i \sim \exp(\beta_j)$$
$$= 1 - \prod_{j=1}^n (1 - e^{-\beta_j t})$$

$$\Rightarrow E(T) = \int_0^\infty t \left[ 1 - \prod_{j=1}^n (1 - e^{-\beta_j t}) \right] dt$$

$$\Rightarrow E(N) = \int_0^\infty t \left[ 1 - \prod_{j=1}^n (1 - e^{-\beta_j t}) \right] dt$$

# Solution

# Birthday Match Problem: Online Version

① People enter the room according to a Poisson process  $N_t$  ( $t \geq 0$ ) with rate  $\lambda = 1$

Each person is independently marked with one of 365 days. w.p.  $\frac{1}{365}$ .

$\Rightarrow$  365 thinned Poisson processes, one for each birthday.

Assume a random person's birthday is uniformly distributed on the 365 days of the year. People enter the room one by one. How many people are in the room the first time that two people share the same birthday? Let  $K$  be the desired number. Find  $E(K)$

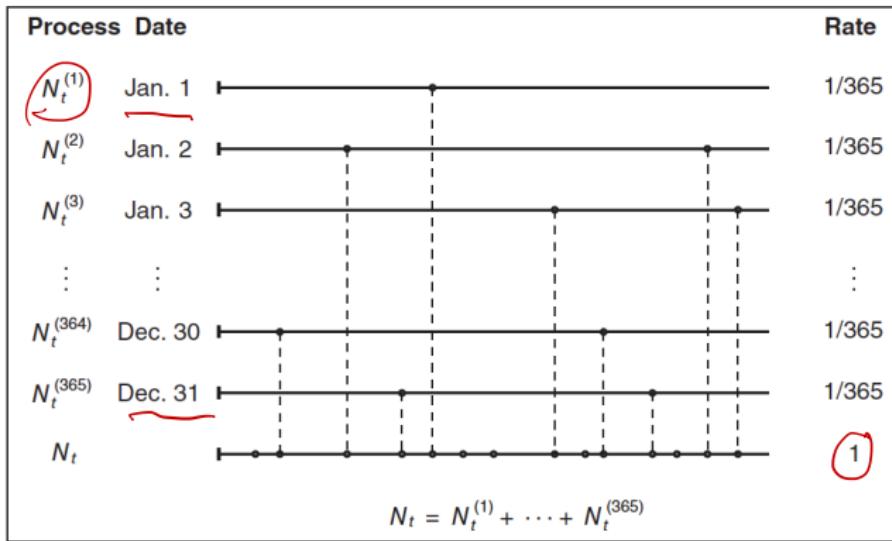
$N_t^{(1)}$  : # of arriving people with birthday at Jan. 1. with rate  $\frac{1}{365}$ .

$N_t^{(365)}$  : - - - - - . . . . . Dec 31 . . . . .  $\frac{1}{365}$ .

# Solution

②  $X_1, \dots, X_n$  : interarrival time of people entering the room.  
 $X_i \sim \text{i.i.d. exp}(1)$ .

Let  $T$  be the first time that two people share the same birthday.



$$T = \sum_{i=1}^k X_i = S_k \quad (k \text{ is the desired number})$$

# Solution

$$T = \sum_{i=1}^k X_i = S_k \quad X_i \text{ independent of } k.$$

$$\Rightarrow E(T) = E\left(\sum_{i=1}^k X_i\right) = E\left[E\left(\sum_{i=1}^k X_i \mid k\right)\right]$$

$$E\left(\sum_{i=1}^k X_i \mid k=n\right) = E\left(\sum_{i=1}^n X_i \mid k=n\right) = E\left(\sum_{i=1}^n X_i\right) = n E(X_i) = n.$$

$$\Rightarrow E\left(\sum_{i=1}^k X_i \mid k\right) = k.$$

$$\Rightarrow E(T) = E[k]$$

(3).  $N_t(\lambda) : S_n \sim \text{Gamma}(n, \lambda).$

for each thinned poisson process  $N_t^{(j)}, j=1, 2, \dots, 365$ .

Let  $Z_j : \text{the first time when the second person marked with birthday } j \text{ enter the room.}$

$$N_{Z_j}^{(1)} = 2$$

$$Z_j = S_2^{(j)} \sim \text{Gamma}\left(2, \frac{1}{365}\right), f_{Z_j}(t) = \frac{t e^{-\frac{t}{365}}}{365^2} \quad t > 0$$

# Solution

CDF of  $Z_1$

$$P(Z_1 \leq t) = \int_0^t \frac{se^{-\frac{s}{365}}}{365^2} ds = 1 - e^{-\frac{t}{365}} \left( 1 + \frac{t}{365} \right)$$

$$\textcircled{4} \quad T = \min(Z_1, \dots, Z_{365})$$

$$P(T > t) = P(\min(Z_1, \dots, Z_{365}) > t) = P(Z_1 > t, \dots, Z_{365} > t)$$

$$= P(Z_1 > t) \cdots P(Z_{365} > t) = [P(Z_1 > t)]^{365} = \left( 1 + \frac{t}{365} \right)^{365} e^{-t}, t > 0$$

$$\Rightarrow E(T) = \int_0^\infty P(T > t) dt = \int_0^\infty \left( 1 + \frac{t}{365} \right)^{365} e^{-t} dt$$

$$\textcircled{5} \quad E(T) = E(K) = \int_0^\infty \left( 1 + \frac{t}{365} \right)^{365} e^{-t} dt = 24.617$$

Thinned poisson process.  
Discrete.  $\leftarrow \frac{\text{poisson}(1)}{\text{bridge.}}$

X

Continuous

Y.

$$E(X) = E(Y) \cdot 1$$

# Solution

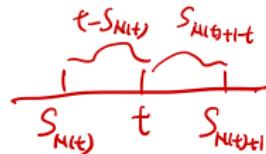
# Feller's Classical Problem

$E(\text{inter-arrival time of two buses}) = 10 \text{ minutes}.$

$E(\text{--- (include time } t\text{)}) \approx 20 \text{ minutes.}$

Buses arrive at a bus stop according to a Poisson process. The time between buses, on average, is 10 minutes. Lisa gets to the bus stop at time  $t$ . How long can she expect to wait for a bus?

# Length-biased Sampling



## Theorem

Suppose buses arrive at a bus stop according to a Poisson process  $N_t$  with parameter  $\lambda$ . Given a fixed  $t > 0$ . The time of the last bus before  $t$  is  $S_{N_t}$ , and the time of the next bus after  $t$  is  $S_{N_t+1}$ . Then we have

$$E(S_{N_t+1} - S_{N_t}) = \frac{2 - e^{-\lambda t}}{\lambda} \stackrel{t \rightarrow \infty}{\approx} \frac{2}{\lambda} = 2(\frac{1}{\lambda}).$$

# Proof

$$N_t \triangleq \underline{N(t)}.$$

1<sup>o</sup>. Compute  $E(S_{N_t+1})$

$$\begin{aligned} E(S_{N(t)+1} | N(t)=k) &= E(S_{k+1} | N(t)=k) && \left( \begin{array}{l} (k+1)^{\text{th}} \text{ arrival time} \\ \text{after time } t \\ \text{independent of } N(t) \end{array} \right) \\ &\stackrel{?}{=} E(S_{k+1}) = \frac{k+1}{\lambda}. \\ \Rightarrow E(S_{N(t)+1} | N(t)) &= \frac{N(t)+1}{\lambda}. \end{aligned}$$

Taking Expectation on both sides,

$$E(S_{N(t)+1}) = E\left(\frac{N(t)+1}{\lambda}\right) = \overbrace{\frac{E(N(t))+1}{\lambda}}^{\rightarrow} = \frac{\lambda t + 1}{\lambda} = t + \frac{1}{\lambda},$$

$$\Rightarrow E(S_{N(t)+1} - t) = \frac{1}{\lambda}, \quad \left( \begin{array}{l} \text{memoryless properties} \\ \text{of exponential distribution} \end{array} \right)$$

# Proof

2°. Compute  $E(S_{N(t)})$

$$E(S_{N(t)} | N(t)=k) = E(S_k | N(t)=k)$$

$$\Rightarrow \forall 0 \leq s \leq t, P(U_{(k)} \leq s) = P(\max(U_1, \dots, U_k) \leq s)$$

$$= P(U_1 \leq s, \dots, U_k \leq s) = P(U_1 \leq s) \dots P(U_k \leq s)$$

$$= [P(U_1 \leq s)]^k = \left(\frac{s}{t}\right)^k.$$

Condition on  $N(t)=k$   
 the  $k$ th arrival time  $S_k \sim U_{(k)}$   
 [the maximum of  $k$  i.i.d.  
 uniform( $t$ )  
 $U_1, \dots, U_k$ ]

$$\Rightarrow E(U_{(k)}) = \int_0^\infty P(U_{(k)} > s) ds = \int_0^t (1 - \frac{s^k}{t^k}) ds = \frac{tk}{k+1}$$

$$\Rightarrow E(S_{N(t)} | N(t)=k) = E(S_k | N(t)=k) = E(U_{(k)}) = \frac{tk}{k+1}$$

$$\Rightarrow E(S_{N(t)} | N(t)) = \frac{t N(t)}{N(t)+1} = t - \frac{t}{N(t)+1}$$

Taking expectation on both sides,

$$E(S_{N(t)}) = E\left(t - \frac{t}{N(t)+1}\right)$$

## Proof

$$\Rightarrow E(S_{N(t)}) = t - t \cdot E(\frac{1}{N(t)+1})$$

$$\begin{aligned} E\left(\frac{1}{N(t)+1}\right) &\stackrel{\text{LOTE}}{=} \sum_{k=0}^{\infty} \frac{1}{k+1} P(N(t)=k) = \sum_{k=0}^{\infty} \frac{1}{k+1} \cdot \frac{e^{-\lambda t} \cdot \lambda^k t^k}{k!} \\ &= \frac{e^{-\lambda t}}{\lambda t} \left[ \sum_{k=1}^{\infty} \frac{(\lambda t)^k}{k!} \right] \\ &= \frac{e^{-\lambda t}}{\lambda t} \left[ \sum_{k=0}^{\infty} \frac{(\lambda t)^k}{k!} - 1 \right] \\ &= \frac{e^{-\lambda t}}{\lambda t} [e^{\lambda t} - 1] = \frac{1 - e^{-\lambda t}}{\lambda t} \end{aligned}$$

$$\begin{aligned} \Rightarrow \text{we have } E(S_{N(t)}) &= t - t \cdot \frac{1 - e^{-\lambda t}}{\lambda t} = t - \frac{1}{\lambda} + \frac{e^{-\lambda t}}{\lambda} \\ \Rightarrow t - E(S_{N(t)}) &= \frac{1 - e^{-\lambda t}}{\lambda}. \end{aligned}$$

# Proof

Thus  $E(S_{N(t)+1} - S_{N(t)})$

$$= (t + \frac{1}{\lambda}) - [t - \frac{1}{\lambda} + \frac{e^{-\lambda t}}{\lambda}]$$

$$= \frac{2 - e^{-\lambda t}}{\lambda} \quad \underline{\lambda t \gg 1}$$

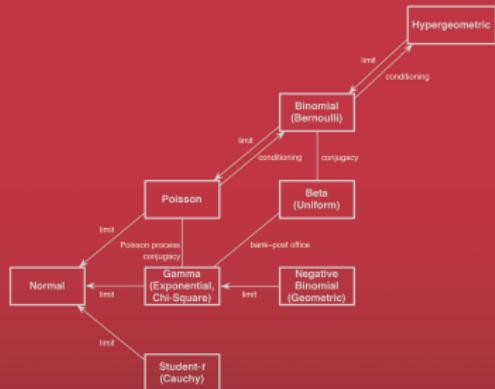
$$\approx \left( \frac{2}{\lambda} \right) \dots$$

## Outline

- 1 Conditional Expectation Given An Event
  - 2 Conditional Expectation Given An R.V.
  - 3 Properties of Conditional Expectation
  - 4 Application I: Prediction
  - 5 Application II: Estimation
  - 6 Application III: Branching Process
  - 7 Application IV: Poisson Process
  - 8 References

Texts in Statistical Science

# Introduction to Probability



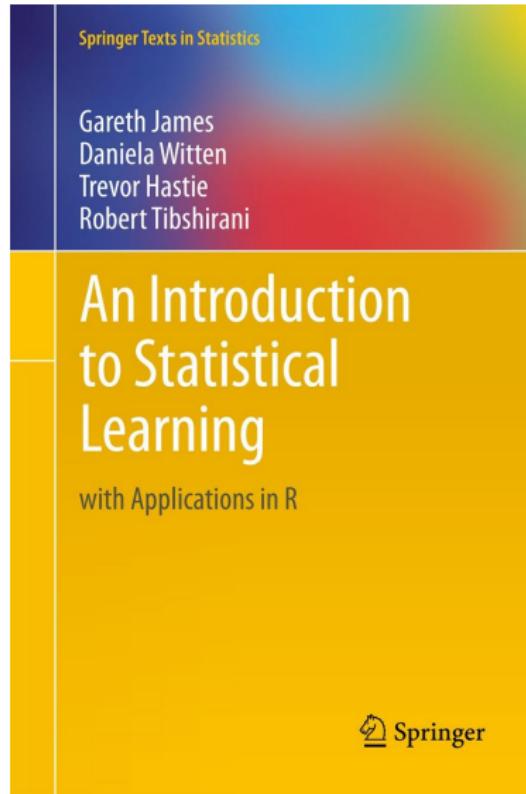
Joseph K. Blitzstein  
Jessica Hwang



CRC Press  
Taylor & Francis Group  
A CHAPMAN & HALL BOOK

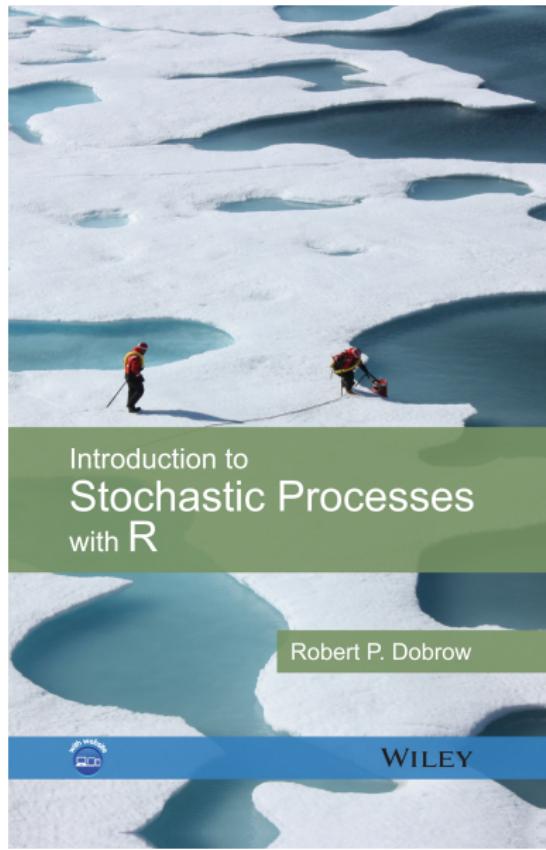
BH

- Introduction to Probability
- Chapman & Hall/CRC, 2014.
- Chapman & Hall/CRC, 2019.
- Chapters 9 & 13



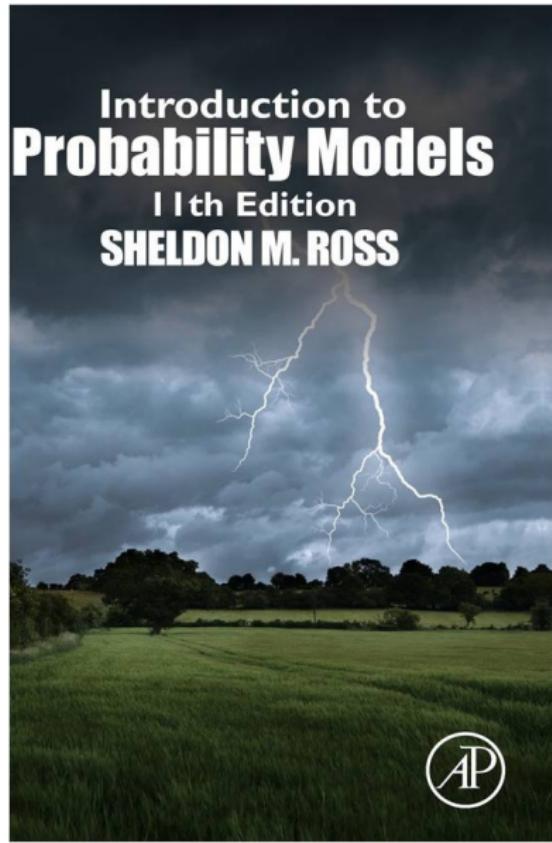
JWHT

- An Introduction to Statistical Learning: with Applications in R
- Springer, 2013.



## SPR

- Introduction to Stochastic Processes with R
- John Wiley & Son, 2016.
- Chapters 1 & 4 & 6



## SMR

- Introduction to Probability Models
- Academic Press, 11 edition, 2014.
- Chapters 1 & 2 & 3 & 5