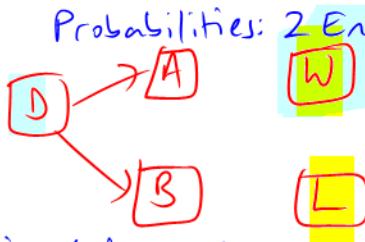


$E(\text{time to win}) = \infty$   
because I might never win.



$P(\text{end up in state } W)$

(not in state L)



44

Expectations, 1 End State (stop)  
interest in  $E(\text{time to stop})$ .

$E(\text{time to FINISH})$  is of interest.

## Chapter 3: Expectation and Variance

In the previous chapter we looked at probability, with three major themes:

1. Conditional probability:  $P(A | B)$ .
2. First-step analysis for calculating eventual probabilities in a stochastic process.
3. Calculating probabilities for continuous and discrete random variables.

In this chapter, we look at the same themes for expectation and variance.

The expectation of a random variable is the long-term average of the random variable  
(e.g. r.v. could be #steps to finish tennis game.)

Imagine observing many thousands of independent random values from the random variable of interest. Take the average of these random values. The expectation is the value of this average as the sample size tends to infinity.

We will repeat the three themes of the previous chapter, but in a different order.

1. Calculating expectations for continuous and discrete random variables.
2. Conditional expectation: the expectation of a random variable  $X$ , *conditional* on the value taken by another random variable  $Y$ . If the value of  $Y$  affects the value of  $X$  (i.e.  $X$  and  $Y$  are *dependent*), the conditional expectation of  $X$  given the value of  $Y$  will be different from the overall expectation of  $X$ .
3. First-step analysis for calculating the expected amount of time needed to reach a particular state in a process (e.g. the expected number of shots before we win a game of tennis).

We will also study similar themes for variance.

$$\frac{3+5+6+6}{4} = 3 + \left(\frac{1}{4}\right) + 5 + \left(\frac{1}{4}\right) + 6 + \left(\frac{2}{4}\right)$$

### 3.1 Expectation

The mean, expected value, or expectation of a random variable  $X$  is written as  $\mathbb{E}(X)$  or  $\mu_X$ . If we observe  $N$  random values of  $X$ , then the mean of the  $N$  values will be approximately equal to  $\mathbb{E}(X)$  for large  $N$ . The expectation is defined differently for continuous and discrete random variables.

*Definition:* Let  $X$  be a continuous random variable with p.d.f.  $f_X(x)$ . The expected value of  $X$  is

$$\rightarrow \mathbb{E}(X) = \int_{-\infty}^{\infty} x f_X(x) dx \quad \text{where } f_X(x) = \underline{\text{PDF of } X}.$$

*Definition:* Let  $X$  be a discrete random variable with probability function  $f_X(x)$ . The expected value of  $X$  is

$$\mathbb{E}(X) = \sum_x x f_X(x) = \sum_x x \underline{P(X=x)} \quad \text{where } f_X(x) \text{ is the PROBABILITY FUNCTION of } X.$$

#### Expectation of $g(X)$

$$\text{eg } g(X) = \sqrt{x}$$

Let  $g(X)$  be a function of  $X$ . We can imagine a long-term average of  $g(X)$  just as we can imagine a long-term average of  $X$ . This average is written as  $\mathbb{E}(g(X))$ . Imagine observing  $X$  many times ( $N$  times) to give results  $x_1, x_2, \dots, x_N$ . Apply the function  $g$  to each of these observations, to give  $g(x_1), \dots, g(x_N)$ . The mean of  $g(x_1), g(x_2), \dots, g(x_N)$  approaches  $\mathbb{E}(g(X))$  as the number of observations  $N$  tends to infinity.

$$\frac{\sqrt{3} + \sqrt{5} + \sqrt{6} + \sqrt{6}}{4}$$

*Definition:* Let  $X$  be a continuous random variable, and let  $g$  be a function. The expected value of  $g(X)$  is

$$\mathbb{E}(g(X)) = \int_{-\infty}^{\infty} g(x) f_X(x) dx. \quad \checkmark$$

*Definition:* Let  $X$  be a discrete random variable, and let  $g$  be a function. The expected value of  $g(X)$  is

$$\mathbb{E}(g(X)) = \sum_x g(x) f_X(x) = \sum_x g(x) \mathbb{P}(X=x).$$

$X = \text{person's height}$   
 $\underline{Y = \text{person's weight}}$

## Expectation of $XY$ : the definition of $\mathbb{E}(XY)$

Suppose we have two random variables,  $X$  and  $Y$ . These might be independent, in which case the value of  $X$  has no effect on the value of  $Y$ . Alternatively,  $X$  and  $Y$  might be *dependent*: when we observe a random value for  $X$ , it might influence the random values of  $Y$  that we are most likely to observe. For example,  $X$  might be the height of a randomly selected person, and  $Y$  might be the weight. On the whole, larger values of  $X$  will be associated with larger values of  $Y$ .

To understand what  $\mathbb{E}(XY)$  means, think of observing a large number of pairs  $(x_1, y_1), (x_2, y_2), \dots, (x_N, y_N)$ . If  $X$  and  $Y$  are dependent, the value  $x_i$  might affect the value  $y_i$ , and vice versa, so we have to keep the observations together in their pairings. As the number of pairs  $N$  tends to infinity, the average  $\frac{1}{N} \sum_{i=1}^N x_i \times y_i$  approaches the expectation  $\mathbb{E}(XY)$ .

For example, if  $X$  is height and  $Y$  is weight,  $\mathbb{E}(XY)$  is the average of (height  $\times$  weight). We are interested in  $\mathbb{E}(XY)$  because it is used for calculating the *covariance* and *correlation*, which are measures of how closely related  $X$  and  $Y$  are (see Section 3.2).

## Properties of Expectation

Revision.

- i) Let  $g$  and  $h$  be functions, and let  $a$  and  $b$  be constants. For any random variable  $X$  (discrete or continuous),

$$\mathbb{E}\{ag(X) + bh(X)\} = a\mathbb{E}\{g(X)\} + b\mathbb{E}\{h(X)\}.$$

In particular,

$$\mathbb{E}(aX + b) = a\mathbb{E}(X) + b.$$

- ii) Let  $X$  and  $Y$  be ANY random variables (discrete, continuous, independent, or non-independent). Then

$$\mathbb{E}(X + Y) = \mathbb{E}(X) + \mathbb{E}(Y).$$

More generally, for ANY random variables  $X_1, \dots, X_n$ ,

$$\mathbb{E}(X_1 + \dots + X_n) = \mathbb{E}(X_1) + \dots + \mathbb{E}(X_n).$$

$X = \#$  cups of tea you had for breakfast.

$Y = \# \dots \text{coffee} \dots \dots \dots$

$$\mathbb{E}X > 0 \quad \mathbb{E}Y > 0 \quad \mathbb{E}(XY) = 0 \neq \mathbb{E}(X)\mathbb{E}(Y).$$

iii) Let  $X$  and  $Y$  be independent random variables, and  $g, h$  be functions. Then

$$\mathbb{E}(XY) = \mathbb{E}(X)\mathbb{E}(Y) \quad \text{ONLY if } X \text{ & } Y \text{ are independent}$$

$$\mathbb{E}\{g(X)h(Y)\} = \mathbb{E}\{g(X)\}\mathbb{E}\{h(Y)\}.$$

**Notes:** 1.  $\mathbb{E}(XY) = \mathbb{E}(X)\mathbb{E}(Y)$  is ONLY generally true if  $X$  and  $Y$  are INDEPENDENT.

2. If  $X$  and  $Y$  are independent, then  $\mathbb{E}(XY) = \mathbb{E}(X)\mathbb{E}(Y)$ . However, the converse is not generally true: it is possible for  $\mathbb{E}(XY) = \mathbb{E}(X)\mathbb{E}(Y)$  even though  $X$  and  $Y$  are dependent.

$$\mathbb{E}X$$

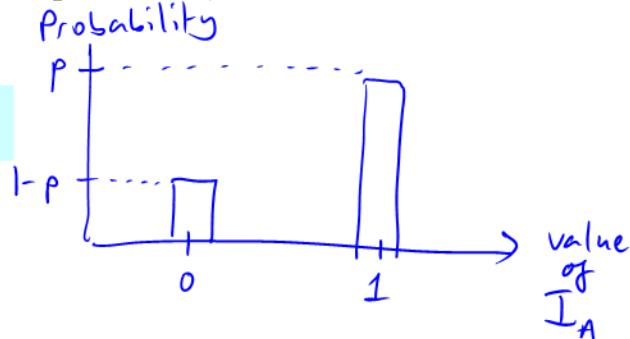
$$\mathbb{P}(A) = \mathbb{E}(I_A)$$

## Probability as an Expectation

Let  $A$  be any event. We can write  $\mathbb{P}(A)$  as an expectation, as follows.

Define the **indicator function**:

$$I_A = \begin{cases} 1 & \text{if event } A \text{ occurs} \\ 0 & \text{otherwise.} \end{cases}$$



Then  $I_A$  is a **random variable**, and

$$\begin{aligned} \mathbb{E}(I_A) &= \sum_{r=0}^1 r \mathbb{P}(I_A = r) \\ &= 0 * \mathbb{P}(I_A = 0) + 1 * \mathbb{P}(I_A = 1) \\ &= \mathbb{P}(I_A = 1) \\ &= \mathbb{P}(A). \end{aligned}$$

Thus

$$\boxed{\mathbb{P}(A) = \mathbb{E}(I_A) \text{ for any event } A.}$$

Office Hours : Rm 208 (upstairs)

Wed : now 12-1pm **every week**

- could not schedule room for 11.30-12.30 : sorry!

Please change in your Coursebooks, p.1. ←

Mon : 2-3pm

Thur : 11.30 - 12.30 pm as before.

### 3.2 Variance, covariance, and correlation

The variance of a random variable  $X$  is a measure of how *spread out* it is. Are the values of  $X$  clustered tightly around their mean, or can we commonly observe values of  $X$  a long way from the mean value? The *variance* measures how far the values of  $X$  are from their mean, on average.

$$\text{Var}(X) = 0 \Leftrightarrow X$$

is a constant.

*Definition:* Let  $X$  be any random variable. The variance of  $X$  is

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

for understanding       $\mathbb{E}(X^2)$       do some algebra      best for using/calculations

The variance is the *mean squared deviation* of a random variable from its own mean.

If  $X$  has *high variance*, we can observe values of  $X$  a long way from the mean.

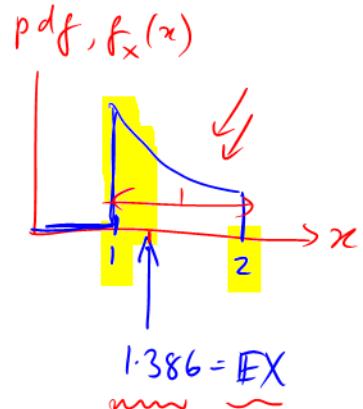
If  $X$  has *low variance*, the values of  $X$  tend to be clustered tightly around the mean value.

*Example:* Let  $X$  be a continuous random variable with p.d.f.

$$f_X(x) = \begin{cases} 2x^{-2} & \text{for } 1 < x < 2, \\ 0 & \text{otherwise.} \end{cases}$$

Find  $\mathbb{E}(X)$  and  $\text{Var}(X)$ .

$$\begin{aligned} \mathbb{E}X &= \int_{-\infty}^{\infty} x f_X(x) dx = \int_1^2 x 2x^{-2} dx \\ &= \int_1^2 2x^{-1} dx \\ &= [2 \log x]_1^2 \\ &= 2 \log 2 - 2 \log 1 \\ &= 2 \log 2 \\ &= \underline{\underline{1.386}}. \end{aligned}$$



$$\log = \log_e = \ln$$

For  $\text{Var}(X)$ , use  $\text{Var}(X) = \mathbb{E}(X^2) - (\mathbb{E}X)^2$

$$\begin{aligned} \text{Now } \mathbb{E}(X^2) &= \int_{-\infty}^{\infty} x^2 f_X(x) dx \\ &= \int_1^2 x^2 2x^{-2} dx \\ &= \int_1^2 2 dx \\ &= [2x]_1^2 \\ &= 2*2 - 2*1 \\ \mathbb{E}(X^2) &= \underline{\underline{2}}. \end{aligned}$$

$$\begin{aligned} \text{So } \text{Var}(X) &= \mathbb{E}(X^2) - (\mathbb{E}X)^2 \\ &= 2 - \{2 \log 2\}^2 \quad \text{from (a)} \\ &= \underline{\underline{0.0782}}. \end{aligned}$$

## Covariance

Covariance is a measure of the association or dependence between two random variables  $X$  and  $Y$ . Covariance can be either positive or negative. (Variance is always positive.)

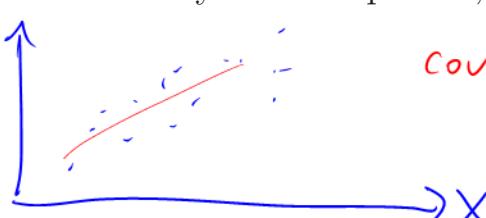
$\hookrightarrow$  (or 0)

Definition: Let  $X$  and  $Y$  be any random variables. The covariance between  $X$  and  $Y$  is given by

$$\text{cov}(X, Y) = \mathbb{E} \left\{ (X - \mu_X)(Y - \mu_Y) \right\} \stackrel{\text{understanding}}{=} \stackrel{\text{algebra}}{=} \stackrel{\text{using}}{=} \mathbb{E}(XY) - (\mathbb{E}X)(\mathbb{E}Y)$$

Where  $\mu_X = \mathbb{E}X$  and  $\mu_Y = \mathbb{E}Y$ .

1.  $\text{cov}(X, Y)$  will be positive if large values of  $X$  tend to occur with large values of  $Y$ , and small values of  $X$  tend to occur with small values of  $Y$ . For example, if  $X$  is height and  $Y$  is weight of a randomly selected person, we would expect  $\text{cov}(X, Y)$  to be positive.



$$\text{cov}(X, Y) > 0$$

e.g.  $X = \text{weight}$ ,  $Y = \text{height}$

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

understanding  $\text{Cov}(X, Y)$

$$= \mathbb{E}\{(X - \mu_X)(Y - \mu_Y)\}$$

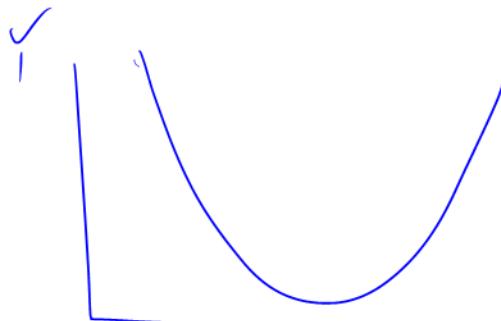
$$\text{Var}(X) = \mathbb{E}(X^2) - (\mathbb{E}X)^2 = \mathbb{E}\{X \cdot X\} - (\mathbb{E}X) \cdot (\mathbb{E}X)$$

algebra

use for calculations

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

If  $X, Y$ , indep  $\mathbb{E}(XY) = (\mathbb{E}X)(\mathbb{E}Y)$  ↗  
 $\Rightarrow \mathbb{E}(XY) - (\mathbb{E}X)(\mathbb{E}Y) = 0$  if indep.



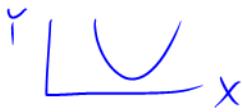
Let  $Y = aX^2 - bX$   
try to find  $a, b$  s.t.

$$\mathbb{E}(XY) = (\mathbb{E}X)(\mathbb{E}Y)$$

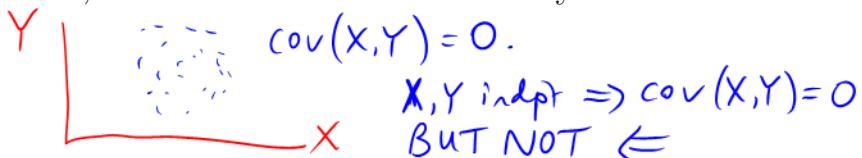
X Possible to construct  
a PERFECT non-linear  
relationship with 0  
covariance.

Warning!  
Covariance only measures LINEAR relationship.

$$\therefore \text{cov}(X, Y) < 0$$



2.  $\text{cov}(X, Y)$  will be **negative** if large values of  $X$  tend to occur with small values of  $Y$ , and small values of  $X$  tend to occur with large values of  $Y$ . For example, if  $X$  is age of a randomly selected person, and  $Y$  is heart rate, we would expect  $X$  and  $Y$  to be negatively correlated (older people have slower heart rates).
3. If  $X$  and  $Y$  are independent, then there is no pattern between large values of  $X$  and large values of  $Y$ , so  $\text{cov}(X, Y) = 0$ . However,  $\text{cov}(X, Y) = 0$  does NOT imply that  $X$  and  $Y$  are independent, unless  $X$  and  $Y$  are Normally distributed.



## Properties of Variance

- i) Let  $g$  be a function, and let  $a$  and  $b$  be constants. For any random variable  $X$  (discrete or continuous),

$$\text{Var}(\underline{a} g(X) + \underline{b}) = a^2 \text{Var}(g(X))$$

In particular,  $\text{Var}(aX + b) = a^2 \text{Var}(X)$ .

- ii) Let  $X$  and  $Y$  be independent random variables. Then

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

MUST BE  
INDEPENDENT.

- iii) If  $X$  and  $Y$  are NOT independent, then

$$\text{Var}(X+Y) = \text{Var}(X) + \text{Var}(Y) + 2\text{cov}(X, Y)$$

## Correlation (non-examinable)

algebra.

$$(x+y)^2 = x^2 + y^2 + 2xy$$

$$\text{Var}(X+Y) = \text{Var}(X) + \text{Var}(Y) + 2\text{cov}(X, Y)$$

The correlation coefficient of  $X$  and  $Y$  is a measure of the linear association between  $X$  and  $Y$ . It is given by the covariance, scaled by the overall variability in  $X$  and  $Y$ . As a result, the correlation coefficient is always between  $-1$  and  $+1$ , so it is easily compared for different quantities.

*Definition:* The correlation between  $X$  and  $Y$ , also called the correlation coefficient, is given by

$$\text{corr}(X, Y) = \frac{\text{cov}(X, Y)}{\sqrt{\text{Var}(X)\text{Var}(Y)}}.$$

The correlation measures linear association between  $X$  and  $Y$ . It takes values only between  $-1$  and  $+1$ , and has the same sign as the covariance.

The correlation is  $\pm 1$  if and only if there is a perfect linear relationship between  $X$  and  $Y$ , i.e.  $\text{corr}(X, Y) = 1 \iff Y = aX + b$  for some constants  $a$  and  $b$ .

The correlation is  $0$  if  $X$  and  $Y$  are independent, but a correlation of  $0$  does not *imply* that  $X$  and  $Y$  are independent.  $\text{not} \iff$

### 3.3 Conditional Expectation and Conditional Variance

Throughout this section, we will assume for simplicity that  $X$  and  $Y$  are discrete random variables. However, exactly the same results hold for continuous random variables too.

$$X = \text{weight} \quad Y = \text{height}. \quad \text{Fix } Y = 170 \text{ cm}$$

Suppose that  $X$  and  $Y$  are discrete random variables, possibly dependent on each other. Suppose that we fix  $Y$  at the value  $y$ . This gives us a set of conditional probabilities  $\mathbb{P}(X = x | Y = y)$  for all possible values  $x$  of  $X$ . This is called the conditional distribution of  $X$ , given that  $Y = y$ .

*Definition:* Let  $X$  and  $Y$  be discrete random variables. The conditional probability function of  $X$ , given that  $Y = y$ , is:

$$\mathbb{P}(X = x | Y = y) = \frac{\mathbb{P}(X = x \text{ AND } Y = y)}{\mathbb{P}(Y = y)}$$

like  
 $\mathbb{P}(A|B) = \frac{\mathbb{P}(A \cap B)}{\mathbb{P}(B)}$

We write the conditional probability function as:

$$f_{X|Y}(x | y) = \mathbb{P}(X = x | Y = y).$$

*Note:* The conditional probabilities  $f_{X|Y}(x | y)$  sum to one, just like any other probability function:

range over  $x \rightarrow$

$$\sum_x \mathbb{P}(X = x | Y = y) = \sum_x \mathbb{P}_{\{Y=y\}}(X = x) = 1,$$

$\uparrow \text{keep } y \text{ fixed}$

using the subscript notation  $\mathbb{P}_{\{Y=y\}}$  of Section 2.3.

We can also find the expectation and variance of  $X$  with respect to this conditional distribution. That is, if we know that the value of  $Y$  is fixed at  $y$ , then we can find the mean value of  $X$  given that  $Y$  takes the value  $y$ , and also the variance of  $X$  given that  $Y = y$ .

**Definition:** Let  $X$  and  $Y$  be discrete random variables. The conditional expectation of  $X$ , given that  $Y = y$ , is

$$\begin{aligned} M_{X|Y=y} &= \mathbb{E}(X | Y=y) = \sum_x x P(X=x | Y=y) \\ &= \sum_x x f_{X|Y}(x | y). \end{aligned}$$

$\mathbb{E}(X | Y = y)$  is the mean value of  $X$ , when  $Y$  is fixed at  $y$ .

e.g.  $X = \text{weight}$ ,  $Y = \text{height}$ ,  $\mathbb{E}(X | Y=166)$  = mean weight FOR PEOPLE of my height; 166cm.

### Conditional expectation as a random variable

The unconditional expectation of  $X$ ,  $\mathbb{E}(X)$ , is just a number,

e.g.  $\mathbb{E}X = 60 \text{ kg}$  or  $\mathbb{E}X = 80 \text{ kg}$ .

The conditional expectation,  $\mathbb{E}(X | Y = y)$ , is a number depending on  $y$ .  
ie. it is a function of  $y$ .

If  $Y$  has an influence on the value of  $X$ , then  $Y$  will have an influence on the average value of  $X$ . So, for example, we would expect  $\mathbb{E}(X | Y = 2)$  to be different from  $\mathbb{E}(X | Y = 3)$ . Or  $\mathbb{E}(X | Y = 140 \text{ cm})$  will be different from  $\mathbb{E}(X | Y = 180 \text{ cm})$ .

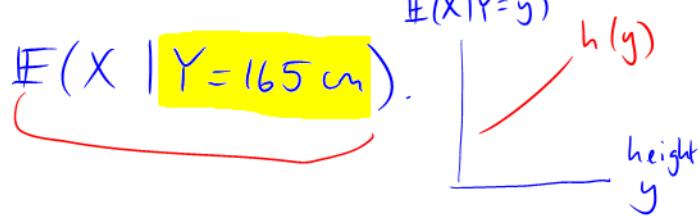
We can therefore view  $\mathbb{E}(X | Y = y)$  as a function of  $y$ , say  $\mathbb{E}(X | Y=y) = h(y)$ .

To evaluate this function,  $h(y) = \mathbb{E}(X | Y = y)$ , we:

i) fix  $Y$  at the chosen value,  $y$ ;

ii) find the average of  $X$ , when  $Y$  is fixed at this value.

e.g. fix  $Y = 165 \text{ cm}$ , find  $\mathbb{E}(X | Y=165 \text{ cm})$ .



$$\mathbb{E}(X) = \sum_x x \boxed{\mathbb{P}(X=x)}$$

"Expectation over the distribution of  $X$ ".

"Expectation over  $X$ "

"Average with respect to  $X$ "

$$\mathbb{E}_x(X) \quad \text{or} \quad \mathbb{E}_x(3XY - 2Y)$$

$$= \sum_x \{ 3xY - 2Y \} \boxed{\mathbb{P}(X=x)}$$

However, we could also evaluate the function at a *random value* of  $Y$ :

- i) observe a **RANDOM** value of  $Y$ ;
- ii) fix  $Y$  at that observed, random value;
- iii) evaluate  $\mathbb{E}(X | Y = \text{observed random value})$ .

We obtain a random variable:  $\mathbb{E}(X | Y) = h(Y)$ .

The randomness comes from the randomness in  $Y$ , not in  $X$ .

Conditional expectation,  $\mathbb{E}(X | Y)$ , is a random variable, with randomness inherited from  $Y$ , not  $X$ .

$\mathbb{E}(X | Y)$  is a convenient but not-completely-intuitive notation so the key thing here is in understanding what the meaning of the notation is.

**Example:** Suppose  $Y = \begin{cases} 1 & \text{with probability } 1/8, \\ 2 & \text{with probability } 7/8, \end{cases}$

and  $X | Y = \begin{cases} 2Y & \text{with probability } 3/4, \\ 3Y & \text{with probability } 1/4. \end{cases}$

Conditional expectation of  $X$  given  $Y = y$  is a number depending on  $y$ :

If  $Y=1$ , then  $X | (Y=1) = \begin{cases} 2 & \text{with prob } 3/4 \\ 3 & \text{.. .. } 1/4 \end{cases}$   
 $\text{so } \mathbb{E}(X | Y=1) = 2 * \frac{3}{4} + 3 * \frac{1}{4} = \frac{9}{4}.$

If  $Y=2$ , then  $X | (Y=2) = \begin{cases} 4 & \text{w.p. } 3/4 \\ 6 & \text{w.p. } 1/4 \end{cases}$   
 $\text{so } \mathbb{E}(X | Y=2) = 4 * \frac{3}{4} + 6 * \frac{1}{4} = \frac{18}{4}.$

So  $\mathbb{E}(X | Y=y) = \begin{cases} 9/4 & \text{if } y=1 \\ 18/4 & \text{if } y=2. \end{cases}$

So  $\mathbb{E}(X | Y=y)$  is a number depending on  $y$ ,  
 ie. a function of  $y$ .

Conditional expectation of  $X$  given random  $Y$  is a random variable.

From above,  $\mathbb{E}(X|Y) = \begin{cases} 9/4 & \text{if } Y=1 \text{ (probability } 1/8) \\ 18/4 & \text{if } Y=2 \text{ (probability } 7/8) \end{cases}$

So  $\underline{\mathbb{E}(X|Y)} = \begin{cases} 9/4 & \text{with prob. } 1/8 \\ 18/4 & \text{with prob. } 7/8. \end{cases}$

These are  $Y$ 's probabilities!

Thus  $\mathbb{E}(X|Y)$  is a random variable.

The randomness (probabilities) come from  $Y$ , not from  $X$ .

Conditional expectation is a very useful tool for finding the *unconditional* expectation of  $X$  (see below). Just like the Partition Theorem, it is useful because it is often easier to specify conditional probabilities than to specify overall probabilities.

## Conditional variance

The conditional variance is similar to the conditional expectation.

- $\text{Var}(X|Y=y)$  is the variance of  $X$ , when  $Y$  is fixed at the value  $Y=y$ .
- $\text{Var}(X|Y)$  is a random variable, giving the variance of  $X$  when  $Y$  is fixed at a value to be selected randomly.

*Definition:* Let  $X$  and  $Y$  be random variables. The conditional variance of  $X$ , given  $Y$ , is given by

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

Like expectation,

$\text{Var}(X|Y=y)$  is a number depending on  $y$  (a function of  $y$ );

$\text{Var}(X|Y)$  is a random variable with randomness inherited from  $Y$ .

Partition Theorem = "Law of Total Probability"

$$P(A) = \sum P(A|B_y) P(B_y)$$

## → Laws of Total Expectation and Variance

If all the expectations below are finite, then for ANY random variables  $X$  and  $Y$ , we have:

i)  $\underline{\mathbb{E}(X)} = \mathbb{E}_Y \{ \underline{\mathbb{E}(X|Y)} \}$  Law of Total Expectation  
LEARN

Note that we can pick ANY r.v.  $Y$ , to make the calculation as easy as possible.

ii)  $\mathbb{E} \{ g(X) \} = \mathbb{E}_Y \{ \mathbb{E}(g(X)|Y) \}$  for any function  $g$ .

iii)  $\mathbb{V}\text{ar}(X) = \mathbb{E}_Y (\mathbb{V}\text{ar}(X|Y)) + \mathbb{V}\text{ar}_Y (\mathbb{E}(X|Y))$

ie.  $\sum_y \mathbb{P}(Y=y)$  Here:  $\sum_y \mathbb{E}(X|Y=y) \mathbb{P}(Y=y) = \mathbb{E}X$

→ Note:  $\mathbb{E}_Y$  and  $\mathbb{V}\text{ar}_Y$  denote expectation over  $Y$  and variance over  $Y$ ,

i.e. the expectation or variance is computed over the distribution of the random variable  $Y$ .

The Law of Total Expectation says that the total average ( $\mathbb{E}X$ ), is the average of case-by-case averages:

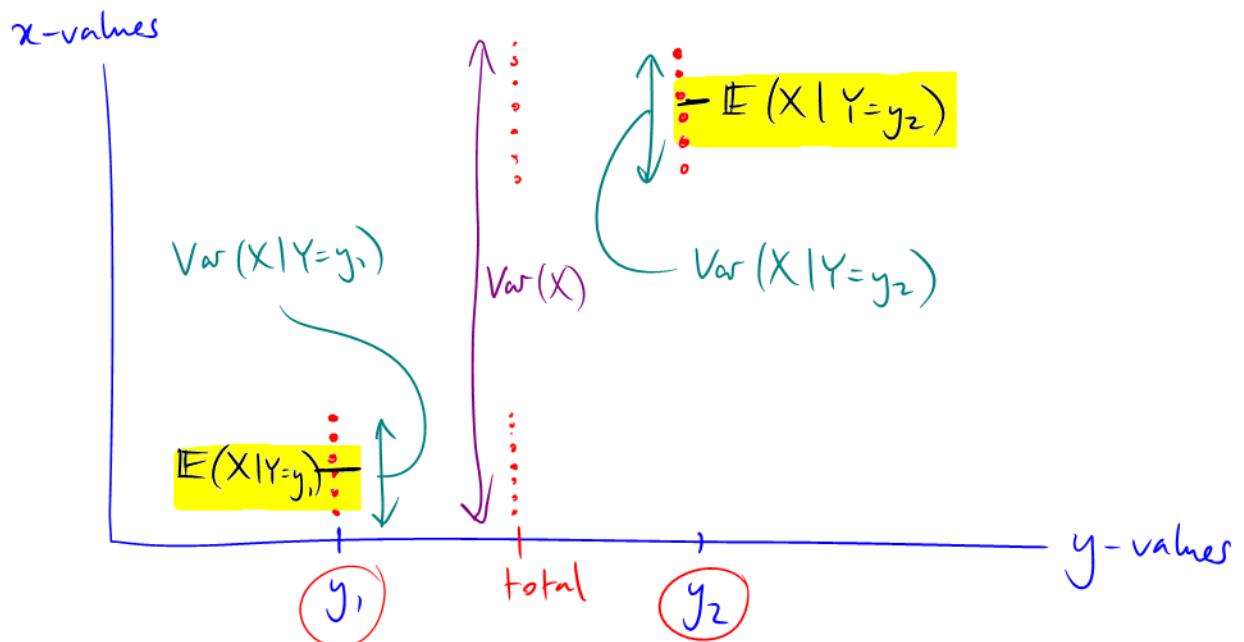
- The total average is  $\mathbb{E}(X)$ ;
- The case-by-case averages are  $\mathbb{E}(X|Y)$  for the different values of  $Y$ ;
- The average of case-by-case averages is average, over the distribution of  $Y$ , of the  $Y$ -case averages;

$$\mathbb{E}_Y \{ \mathbb{E}(X|Y) \} = \sum_y \mathbb{E}(X|Y=y) \mathbb{P}(Y=y) = \mathbb{E}(h(Y))$$

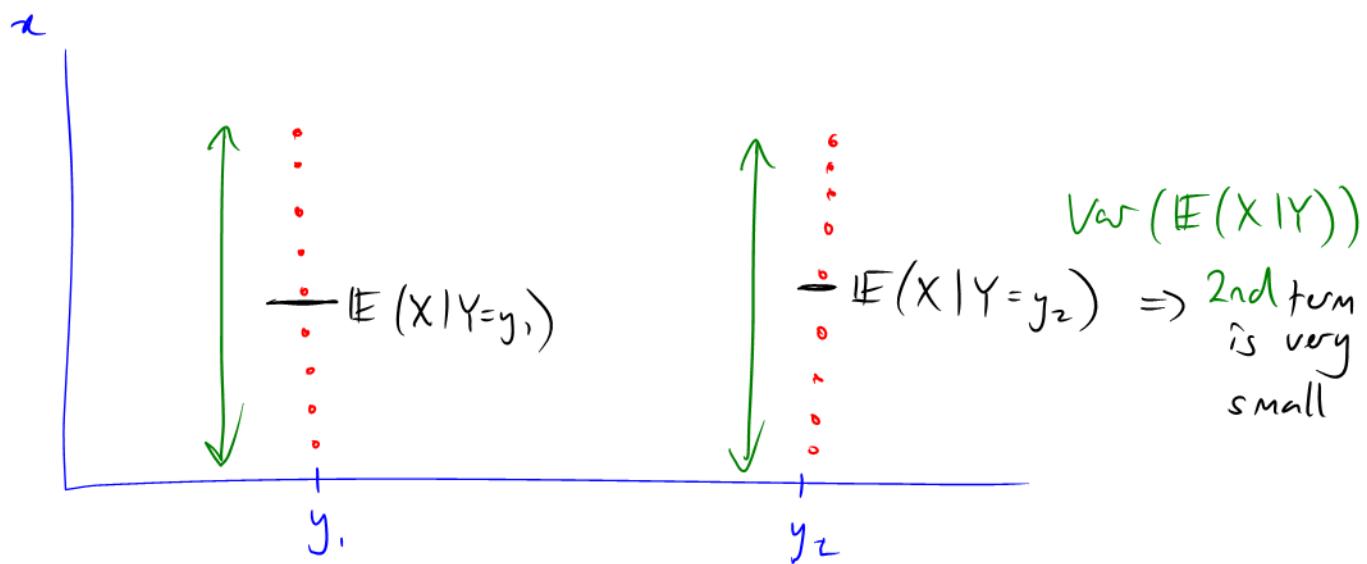
NOTATION: (moderately standard):  $\mathbb{E}_Y$  = average with respect to the distribution of  $Y$ .

Today & Tomorrow;  
no office now, sorry!

Use email.



$$Var(X) = \underbrace{E_Y(Var(X|Y))}_{\text{average within-group variance}} + \underbrace{Var_Y(E(X|Y))}_{\text{between-group variance}} \quad .4$$



1st term (within group)  
will dominate in this example.

$$\text{total Expectation} \hookrightarrow \mathbb{E}(X) = \mathbb{E}_Y \{ \mathbb{E}(X|Y) \}$$

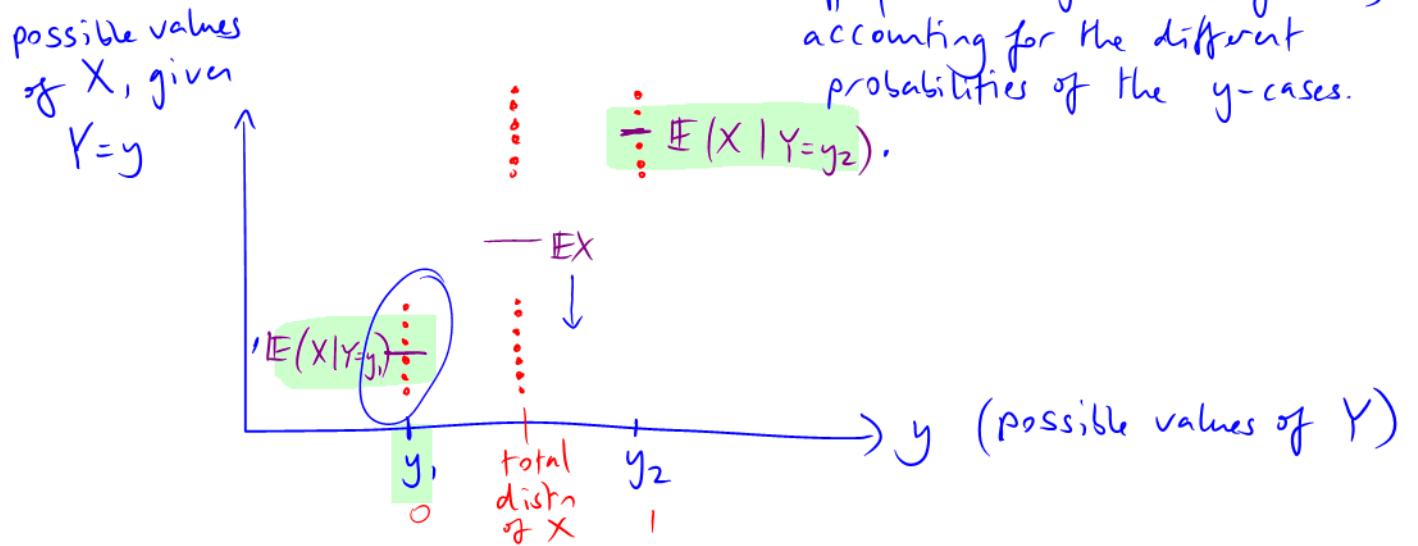
expectation with respect to  $Y$

conditional expectation

$$= \left( \sum_y \mathbb{E}(X|Y=y) P(Y=y) \right)$$

case-by-case expectations of  $X$  for particular  $y$ -values

take an appropriate weighted sum of these, accounting for the different probabilities of the  $y$ -cases.



$Y$  continuous       $X$  = discrete

$$\begin{aligned} \mathbb{P}(X=x) &= \int_y \mathbb{P}(X=x | Y=y) f_Y(y) dy \\ &= \mathbb{E}_Y \{ \mathbb{P}(X|Y) \} \end{aligned}$$

**Example:** In the example above, we had:  $\mathbb{E}(X | Y) = \begin{cases} 9/4 & \text{with probability } 1/8, \\ 18/4 & \text{with probability } 7/8. \end{cases}$

The total average is:

$$\mathbb{E}(X) = \mathbb{E}_Y \{ \mathbb{E}(X | Y) \} = \frac{9}{4} * \frac{1}{8} + \frac{18}{4} * \frac{7}{8} = 4.22.$$


---

**Proof of (i), (ii), (iii):**

(i) is a special case of (ii), so we just need to prove (ii). Begin at RHS:

$$\begin{aligned}
 \text{RHS} &= \mathbb{E}_Y [\mathbb{E}(g(X) | Y)] = \mathbb{E}_Y \left[ \sum_x g(x) \mathbb{P}(X = x | Y) \right] \\
 &= \sum_y \left[ \sum_x g(x) \mathbb{P}(X = x | Y = y) \right] \mathbb{P}(Y = y) \\
 &= \sum_y \sum_x g(x) \mathbb{P}(X = x | Y = y) \mathbb{P}(Y = y) \\
 &= \sum_x g(x) \sum_y \mathbb{P}(X = x | Y = y) \mathbb{P}(Y = y) \\
 &= \sum_x g(x) \mathbb{P}(X = x) \quad (\text{partition rule}) \\
 &= \mathbb{E}(g(X)) = \text{LHS}.
 \end{aligned}$$

(iii) Wish to prove  $\text{Var}(X) = \mathbb{E}_Y[\text{Var}(X | Y)] + \text{Var}_Y[\mathbb{E}(X | Y)]$ . Begin at RHS:

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

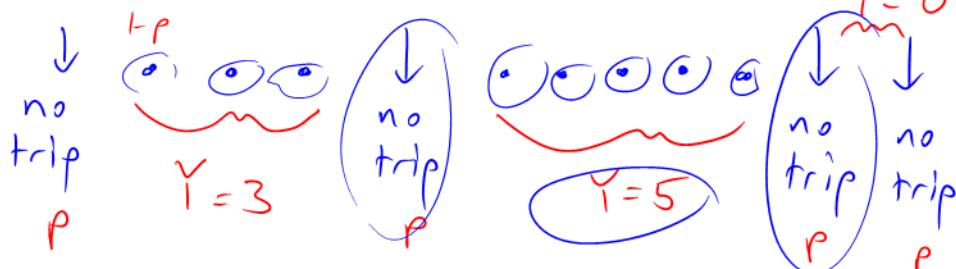

---

$$= \mathbb{E}_Y \left\{ \mathbb{E}(X^2 | Y) - (\mathbb{E}(X | Y))^2 \right\} + \left\{ \mathbb{E}_Y \left\{ [\mathbb{E}(X | Y)]^2 \right\} - \left[ \underbrace{\mathbb{E}_Y(\mathbb{E}(X | Y))}_{\mathbb{E}(X) \text{ by part (i)}} \right]^2 \right\}$$

$$= \underbrace{\mathbb{E}_Y \{ \mathbb{E}(X^2 | Y) \}}_{\mathbb{E}(X^2) \text{ by part (i)}} - \mathbb{E}_Y \{ [\mathbb{E}(X | Y)]^2 \} + \mathbb{E}_Y \{ [\mathbb{E}(X | Y)]^2 \} - (\mathbb{E}X)^2$$

$$= \mathbb{E}(X^2) - (\mathbb{E}X)^2$$

$$= \text{Var}(X) = \text{LHS}. \quad \square$$



### 3.4 Examples of Conditional Expectation and Variance

#### 1. Swimming with dolphins

Fraser runs a dolphin-watch business.

Every day, he is unable to run the trip

due to bad weather with probability  $p$ ,

independently of all other days. Fraser works every day except the bad-weather days, which he takes as holiday.



Let  $Y$  be the number of consecutive days Fraser has to work between bad-weather days. Let  $X$  be the total number of customers who go on Fraser's trip in this period of  $Y$  days. Conditional on  $Y$ , the distribution of  $X$  is

$$(X | Y) \sim \text{Poisson}(\mu Y).$$

**Q1 Ass)** (a) Name the distribution of  $Y$ , and state  $\mathbb{E}(Y)$  and  $\text{Var}(Y)$ .

(b) Find the expectation and the variance of the number of customers Fraser sees between bad-weather days,  $\mathbb{E}(X)$  and  $\text{Var}(X)$ .

(a) Let "success" = "bad-weather day" ↙  
 "failure" = "work-day"

$Y = \# \text{failures before first success}$

So  $Y \sim \text{Geometric}(p)$

$$\begin{aligned} \mathbb{P}(\text{success}) &= \mathbb{P}(\text{bad weather}) \\ &= p. \end{aligned}$$

$$\begin{aligned} \text{So } \mathbb{E}(Y) &= \left( \frac{1-p}{p} \right) \\ \text{Var}(Y) &= \frac{1-p}{p^2} \end{aligned} \quad \left. \begin{array}{l} \text{either memory} \\ \text{or Stats 325 exam formula sheet!} \end{array} \right\}$$

(b) We know  $(X | Y) \sim \text{Poisson}(\mu Y)$

$$\text{So } \mathbb{E}(X | Y) = \mu Y$$

$$\text{and } \text{Var}(X | Y) = \mu Y$$

} remember  $\mathbb{E}(X | Y)$  is a r.v.  
 which is a function of  $Y$   
 (same for  $\text{Var}(X | Y)$ ).

## Law of Total Expectation:

$$\begin{aligned}
 E(X) &= E_Y \{ E(X|Y) \} \\
 &= E_Y \{ \mu_Y \} \\
 &= E(\mu_Y) \\
 &= \mu E Y \\
 E(X) &= \frac{\mu(1-p)}{p}
 \end{aligned}$$

quietly dropping the subscript  $Y$ ,  
ie.  $E_Y \rightarrow E$ , because I  
don't need it any more: there's only  
 $Y$  left on the RHS, so I no longer  
need the clarification  $E_Y$ .

using  $EY = \frac{1-p}{p}$  from (a).

## Law of Total Variance:

$$\begin{aligned}
 \text{Var}(X) &= E_Y \left\{ \underbrace{\text{Var}(X|Y)}_{\mu_Y} \right\} + \underbrace{\text{Var}_Y (\underbrace{E(X|Y)}_{\mu_Y})}_{\text{from (a)}} \\
 &= E_Y (\mu_Y) + \text{Var}_Y (\mu_Y) \\
 &= \mu E Y + \mu^2 \text{Var} Y \\
 &= \mu \left( \frac{1-p}{p} \right) + \mu^2 \left( \frac{1-p}{p^2} \right) \text{ from (a)} \\
 &= \frac{\mu(1-p)(p+\mu)}{p^2} = \text{Var}(X). \quad \square
 \end{aligned}$$

variance due to  
different #'s days,  
ie. between-group var.  
it will dominate as  
long as  $\frac{\mu}{p} > 1$

### Checking your answer in R:

If you know how to use a statistical package like  $R$ , you can check your answer to the question above as follows.

```

> # Pick a value for p, e.g. p = 0.2.
> # Pick a value for mu, e.g. mu = 25
>
Read
> # Generate 10,000 random values of Y ~ Geometric(p = 0.2):
> y <- rgeom(10000, prob=0.2)
>
> # Generate 10,000 random values of X conditional on Y:
> # use (X | Y) ~ Poisson(mu * Y) ~ Poisson(25 * Y)
> x <- rpois(10000, lambda = 25*y)

```

```

> # Find the sample mean of X (should be close to E(X)):
> mean(x)
[1] 100.6606
>
> # Find the sample variance of X (should be close to var(X)):
> var(x)
[1] 12624.47
>
> # Check the formula for E(X):
> 25 * (1 - 0.2) / 0.2
[1] 100
>
> # Check the formula for var(X):
> 25 * (1 - 0.2) * (0.2 + 25) / 0.2^2
[1] 12600
  
```

The formulas we obtained by working give  $\mathbb{E}(X) = 100$  and  $\text{Var}(X) = 12600$ . The sample mean was  $\bar{x} = 100.6606$  (close to 100), and the sample variance was 12624.47 (close to 12600). Thus our working seems to have been correct.

## 2. Randomly stopped sum

This model arises very commonly in stochastic processes. A random number  $N$  of events occur, and each event  $i$  has associated with it some cost, penalty, or reward  $X_i$ . The question is to find the mean and variance of the total cost / reward:

$T_N = X_1 + X_2 + \dots + X_N$ .  
e.g.  $X_i = \text{amount of money withdrawn by customer } i$

The difficulty is that the number  $N$  of terms in the sum is itself random.



$T_N$  is called a randomly stopped sum: it is a sum of  $X_i$ 's, randomly stopped at the random number of terms,  $N$ .

**Example:** Think of a cash machine, which has to be loaded with enough money to cover the day's business. The number of customers per day is a random number  $N$ . Customer  $i$  withdraws a random amount  $X_i$ . The total amount withdrawn during the day is a randomly stopped sum:  $T_N = X_1 + \dots + X_N$ .

## Cash machine example

The citizens of Remuera withdraw money from a cash machine according to the following probability function ( $X$ ):

Amount, $x$ (\$)	50	100	200
$\mathbb{P}(X = x)$	0.3	0.5	0.2

So  $\mathbb{E}N = \lambda$

$\checkmark \text{Var}(N) = \lambda$

The number of customers per day has the distribution  $N \sim \text{Poisson}(\lambda)$ .

Let  $T_N = X_1 + X_2 + \dots + X_N$  be the total amount of money withdrawn in a day, where each  $X_i$  has the probability function above, and  $X_1, X_2, \dots$  are independent of each other and of  $N$ .

$T_N$  is a randomly stopped sum, stopped by the random number of  $N$  customers.

(a) Show that  $\mathbb{E}(X) = 105$ , and  $\text{Var}(X) = 2725$ .

→ (b) Find  $\mathbb{E}(T_N)$  and  $\text{Var}(T_N)$ : the mean and variance of the amount of money withdrawn each day.

## Solution

(a) Exercise.

(b) Let  $T_N = \sum_{i=1}^N X_i$ . If we knew how many terms,  $N$ , were in the sum, it would be easy to find  $\mathbb{E}(T_N)$  and  $\text{Var}(T_N)$  as the mean & variance of a sum of independent r.v.s,  $X_1, \dots, X_N$ . So "pretend" we know how many terms there are, i.e. condition on  $N$ .

condition on  $N$  on LHS

$\mathbb{E}(T_N | N) \Rightarrow N$  is like a CONSTANT on the RHS

$$\rightarrow \mathbb{E}(T_N | N) = \mathbb{E}\{X_1 + X_2 + \dots + X_N | N\}$$

$$= \mathbb{E}\{X_1 + X_2 + \dots + X_N\} \quad \text{because all } X_i \text{'s are independent of } N$$

emphasises that, once we know  $N = \# \text{terms}$ , there is no further influence of  $N$  on the distn of  $X_i$ 's, e.g.  $\mathbb{E}(X_i | N) = \mathbb{E}(X_i)$

Do NOT need  
indep of  $X_i$ 's for  
this

$$= \mathbb{E}(X_1) + \mathbb{E}(X_2) + \dots + \mathbb{E}(X_N) \quad \text{where } N \text{ is now considered constant}$$

$$= N * \mathbb{E}(X) \quad \text{because all } X_i \text{'s have same mean, } \mathbb{E}X$$

$$\mathbb{E}(T_N | N) = 105N. \quad \text{by part (a)} \quad \mathbb{E}X = 105.$$

Similarly,

$$\text{Var}(T_N | N) = \text{Var}(X_1 + \dots + X_N | N)$$

$$= \text{Var}(X_1 + \dots + X_N) \text{ because } X_i \text{'s are indept of } N$$

$$= \text{Var}(X_1) + \text{Var}(X_2) + \dots + \text{Var}(X_N) \text{ because } X_i \text{'s indept of each other.}$$

(DO NEED INDEP of  $X_i$ 's for this)

$$= N \text{Var}(X) \quad - \text{all } X_i \text{'s have same variance, } \text{Var}(X)$$

$$\therefore \text{Var}(T_N | N) = 2725N$$

$$\text{So } \mathbb{E}(T_N) = \mathbb{E}_N \{ \mathbb{E}(T_N | N) \}$$

$$= \mathbb{E}_N \{ 105N \} \quad \text{from above}$$

$$= 105 \mathbb{E}_N(N) \leftarrow$$

$$\therefore \mathbb{E}(T_N) = 105\lambda \quad \text{because } N \sim \text{Poisson}(\lambda) \text{ so } \mathbb{E}N = \lambda.$$

$$\text{Similarly, } \text{Var}(T_N) = \mathbb{E}_N \{ \text{Var}(T_N | N) \} + \text{Var}_N \{ \mathbb{E}(T_N | N) \}$$

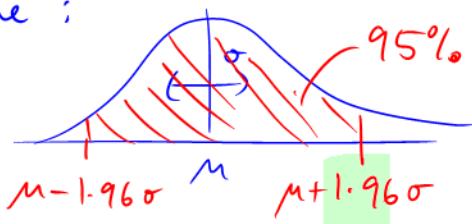
$$= \mathbb{E}_N \{ 2725N \} + \text{Var}_N \{ 105N \}$$

$$= 2725 \mathbb{E}_N(N) + (105)^2 \text{Var}(N)$$

$$= 2725\lambda + 11025\lambda \quad \text{because } N \sim \text{Poisson}(\lambda) \\ \text{so } \mathbb{E}N = \text{Var}(N) = \lambda.$$

$$\therefore \text{Var}(T_N) = 13750\lambda$$

For a quick way to convert from knowing  $\mathbb{E}(T_N)$  and  $\text{Var}(T_N)$  to knowing how much money to put in the machine:



(very rough but probably OK approx due to CLT)

$$\text{use } \mathbb{E}(T_N) + 1.96 \sqrt{\text{Var}(T_N)}$$

and this upper value should cover about 97.5% of scenarios (ie. of days' business).

e.g. If one customer every 5 mins for 8 hrs/day  $\Rightarrow \lambda = 96 \Rightarrow \$12,331 \text{ per day}$

This term is MUCH bigger; it's  $\text{Var}_N$  (total for a fixed  $N$ )  
ie, the variability due to the random #customers is much more important than the variability you'd see for a fixed #customers.

## Check in R (advanced)

```
> # Create a function tn.func to calculate a single value of T_N
> # for a given value N=n:
> tn.func <- function(n){
  sum(sample(c(50, 100, 200), n, replace=T,
  prob=c(0.3, 0.5, 0.2)))
}

> # Generate 10,000 random values of N, using lambda=50:
> N <- rpois(10000, lambda=50)
> # Generate 10,000 random values of T_N, conditional on N:
> TN <- sapply(N, tn.func)
> # Find the sample mean of T_N values, which should be close to
> # 105 * 50 = 5250:
> mean(TN)
[1] 5253.255
> # Find the sample variance of T_N values, which should be close
> # to 13750 * 50 = 687500:
> var(TN)
[1] 682469.4
```

All seems well. Note that the sample variance is often some distance from the true variance, even when the sample size is 10,000.

### General result for randomly stopped sums:

Suppose  $X_1, X_2, \dots$  each have the same mean  $\mu$  and variance  $\sigma^2$ , and  $X_1, X_2, \dots$ , and  $N$  are mutually independent. Let  $T_N = X_1 + \dots + X_N$  be the randomly stopped sum. By following similar working to that above:

$$\mathbb{E}(T_N) = \mathbb{E} \left\{ \sum_{i=1}^N X_i \right\} = \mu \mathbb{E}(N)$$

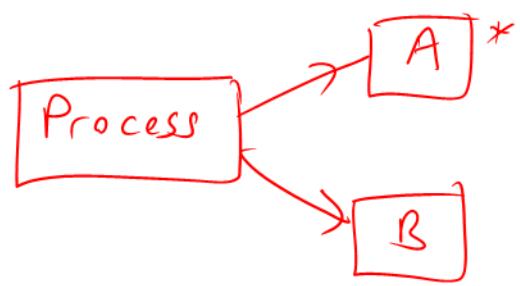
$$\text{Var}(T_N) = \text{Var} \left\{ \sum_{i=1}^N X_i \right\} = \sigma^2 \mathbb{E}(N) + \mu^2 \text{Var}(N).$$

721 students: a simple example of Wald's Equation

Component due to variance in  $X_i$ 's  
 $(\text{Var}(X_i)=\sigma^2)$ :  
 within- $N$  component

component due to variance between different  $N$  values.

Ch 2



of interest:

$P(\text{eventually reach } A)$

i.e.  $P(\text{end up at } A)$

FSA.

If I asked  $\mathbb{E}(\text{time to reach } A) = \infty$

Ch 3



$\mathbb{E}(\text{time to finish})$

- assuming we definitely do finish, otherwise  $\mathbb{E}(\text{time}) = \infty$ .

b/c  $\mathbb{E}(\text{time}) = \infty * p_0 + \dots = \infty$

### 3.5 First-Step Analysis for calculating expected reaching times

Remember from Section 2.6 that we use First-Step Analysis for finding the probability of eventually reaching a particular state in a stochastic process. First-step analysis for probabilities uses **conditional probability** and **The Partition Theorem (= Law of Total Probability)**.

In the same way, we can use first-step analysis for finding the **expected reaching time for a state**.

This is the expected number of steps that will be needed to reach a particular state from a specified start-point, or the expected length of time it will take to get there if we have a continuous time process.

Just as first-step analysis for probabilities uses conditional probability and the law of total probability (Partition Theorem), first-step analysis for expectations uses **conditional expectation** and the **Law of Total Expectation**.

#### First-step analysis for probabilities:

The first-step analysis procedure for probabilities can be summarized as follows:

$$P(\text{eventual goal}) = \sum_{\substack{\text{first-step} \\ \text{options}}} P(\text{eventual goal} | \text{1st step option}) P(\text{1st step option})$$

↑

This is because the first-step options form a **partition of the sample space**,  
 $\Omega = \{ \text{all paths through the system from start-point to end} \}$ .

#### First-step analysis for expected reaching times:

The expression for expected reaching times is very similar:

$$\begin{aligned} \mathbb{E}\{g(X)\} &= \sum_x g(x) P(X=x) \\ \downarrow & \\ \mathbb{E}(\text{reaching time}) &= \sum_{\substack{\text{first-step} \\ \text{options}}} \mathbb{E}(\text{reaching time} | \text{1st step option}) P(\text{1st step option}) \\ &= \mathbb{E}_{\text{1st step options}} \{ \mathbb{E}(\text{reaching time} | \text{option}) \}. \end{aligned}$$

This follows immediately from the law of total expectation:

$$\begin{aligned} \mathbb{E}(X) &= \mathbb{E}_Y \{ \mathbb{E}(X|Y) \} = \mathbb{E} \{ h(Y) \} \\ &= \sum_y h(y) P(Y=y) \\ &= \sum_y \mathbb{E}(X|Y=y) P(Y=y). \end{aligned}$$

Let  $X$  be the reaching time, and let  $Y$  be the label for possible options:

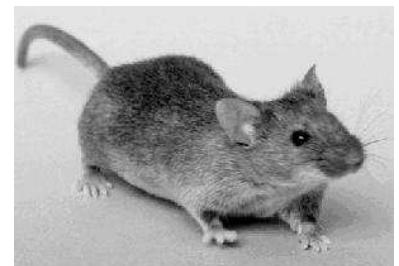
e.g.  $Y = 1, 2, 3, \dots$  for first-step options 1, 2, 3, ... .

We then obtain:

$$\begin{aligned} \mathbb{E}(\text{reaching time}) &= \mathbb{E}(X) \\ &= \sum_y \mathbb{E}(X|Y=y) P(Y=y) \\ &= \sum_{\substack{\text{1st step} \\ \text{options}}} \mathbb{E}(\text{reaching time} \mid \text{1st step option}) P(\text{1st step option}). \end{aligned}$$


---

### Example 1: Mouse in a Maze



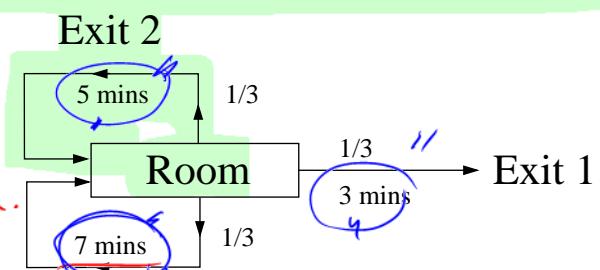
A mouse is trapped in a room with three exits at the centre of a maze.

- Exit 1 leads outside the maze after 3 minutes.
- Exit 2 leads back to the room after 5 minutes.
- Exit 3 leads back to the room after 7 minutes.

Every time the mouse makes a choice, it is equally likely to choose any of the three exits. What is the expected time taken for the mouse to leave the maze?

Let

$X$  = time taken for mouse to leave maze, starting from Room, R.



Let  $Y$  = exit the mouse chooses first (1, 2, or 3).

$$\text{Then: } \mathbb{E}(X) = \mathbb{E}_Y \{ \mathbb{E}(X|Y) \} \quad \text{L o T E}$$

$$= \sum_{y=1}^3 \mathbb{E}(X|Y=y) \mathbb{P}(Y=y)$$

$$= \mathbb{E}(X|Y=1) * \frac{1}{3} + \mathbb{E}(X|Y=2) * \frac{1}{3} + \mathbb{E}(X|Y=3) * \frac{1}{3}$$

But:  $\mathbb{E}(X|Y=1) = 3 \text{ mins}$

$$\mathbb{E}(X|Y=2) = 5 + \mathbb{E}(X)$$

$$\mathbb{E}(X|Y=3) = 7 + \mathbb{E}(X)$$

So:

$$\begin{aligned} \mathbb{E}(X) &= 3 * \frac{1}{3} + \{ 5 + \mathbb{E}(X) \} * \frac{1}{3} \\ &\quad + \{ 7 + \mathbb{E}(X) \} * \frac{1}{3} \end{aligned}$$

$$= 15 * \frac{1}{3} + 2 \mathbb{E}(X) * \frac{1}{3}$$

$$\Rightarrow \frac{1}{3} \mathbb{E}(X) = 15 + \frac{1}{3}$$

$$\Rightarrow \underline{\underline{\mathbb{E}(X) = 15 \text{ minutes}}}$$

→ after 5 mins, get back to Room, from which point we have STILL time  $X_1$  to get out, where  $X_1 \sim X$  so  $\mathbb{E}(X_1) = \mathbb{E}(X)$   
 - note that of course  $X_1 \neq X$ , they have different VALUES!  
 In fact,  $\mathbb{E}(X) = \mathbb{E}(X_1) + 5$   
 $X|Y=2 = X_1 + 5$  where  $X_1 \sim X$

**Important! Use this Approach**

**Notation for quick solutions of first-step analysis problems**

As for probabilities, first-step analysis for expectations relies on a good notation.

The best way to tackle the problem above is as follows.

Define  $m_R = \mathbb{E}(\text{time to leave maze} \mid \text{start in state Room})$

First-step analysis:

$$m_R = \frac{1}{3} * 3 + \frac{1}{3} (5 + m_R) + \frac{1}{3} (7 + m_R)$$

$$\Rightarrow 3m_R = (3 + 5 + 7) + 2m_R$$

$$\Rightarrow \underline{\underline{m_R = 15 \text{ minutes as before.}}}$$

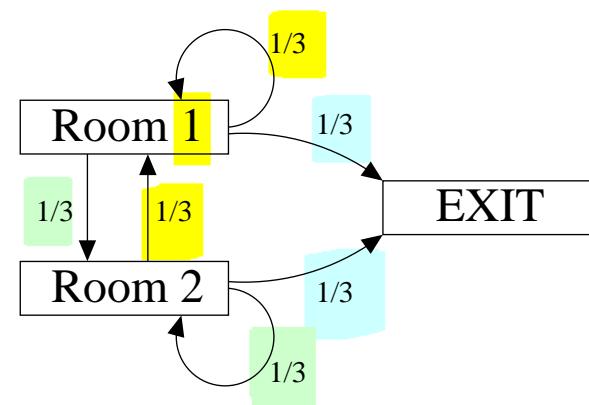
## Example 2: Counting the steps

The most common questions involving first-step analysis for expectations ask for the **expected number of steps before finishing**. The number of steps is usually equal to the **number of arrows traversed from the current state to the end**.

The key point to remember is that when we take expectations, we are usually **counting something**.

You must remember to **add on whatever you are counting, to every step taken.**

The mouse is put in a new maze with two rooms, pictured here. Starting from Room 1, what is the expected number of steps the mouse takes before it reaches the exit? **#arrows**



### 1. Define notation:

$$\text{let } M_1 = \mathbb{E}(\# \text{steps taken to finish} \mid \text{start in Room 1})$$

$$M_2 = \mathbb{E}(\# \text{steps taken to finish} \mid \text{start in Room 2})$$

**Use "m" for "Mean" or  $\mathbb{E}$**

### 2. FSA:

$$m_1 = \frac{1}{3} * (1 + m_1) + \frac{1}{3} * (1 + m_2) + \frac{1}{3} * 1 \quad (a)$$

$$m_2 = \frac{1}{3} * (1 + m_1) + \frac{1}{3} * (1 + m_2) + \frac{1}{3} * 1 \quad (b)$$

Generally solve simultaneously, but here we can see RHS (a) is identical to RHS (b),

$$\text{so } m_1 = m_2 .$$

$$\begin{aligned} \text{Subst in (a)} \Rightarrow m_1 &= \frac{1}{3} (1 + m_1) + \frac{1}{3} (1 + m_1) + \frac{1}{3} \\ &\Rightarrow m_1 = 3 \text{ steps.} \end{aligned}$$

$$\text{So } M_1 = \mathbb{E}(\# \text{steps} \mid \text{start @ 1}) = 3$$

$$\& M_2 = \mathbb{E}(\# \text{steps} \mid \text{start @ 2}) = M_1 = 3 \text{ steps also.}$$

Say  $Y$  is a discrete r.v.  $A$  is any event

Partition Thm  $\Rightarrow P(A) = \sum_y P(A|Y=y) P(Y=y)$

Does that mean that if  $Y$  is continuous,

then  $P(A) = \int_y P(A|Y=y) f_Y(y) dy$

$= \mathbb{E}_Y \{ P(A|Y) \}$  pdf of  $Y$

$$M_x = \mathbb{E}(\text{#steps to exit maze} \mid \text{start in Room } x) \quad x=1, 2$$

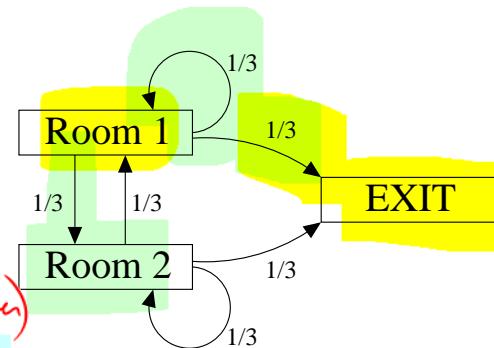
↑ counting arrows

### Incrementing before partitioning

In many problems, all possible first-step options incur the same initial penalty.

The last example is such a case, because  
 every possible step adds 1 arrow  
 to the total #steps taken (#arrows taken)

In a case where all steps incur the same penalty, there are two ways of proceeding:



→ 1. Add the penalty (add each arrow as it occurs) onto each of the 1st-step options separately:

$$\text{e.g. } M_1 = \frac{1}{3} * 1 + \frac{1}{3} (1 + M_2) + \frac{1}{3} (1 + M_1)$$

2. (Usually quicker): add the penalty (e.g. 1 arrow here) once only, at the beginning:

*not always possible e.g. Ass2 Q1c*

$$M_1 = 1 + \frac{1}{3} * 0 + \frac{1}{3} * M_2 + \frac{1}{3} M_1$$

In each case, we will get the same answer (check). This is because the option probabilities sum to 1, so in Method 1 we are adding  $(\frac{1}{3} + \frac{1}{3} + \frac{1}{3}) * 1 = 1 * 1 = 1$ , which is the same as we are adding in Method 2.

### 3.6 Probability as a conditional expectation

Recall from Section 3.1 that for any event  $A$ , we can write  $\mathbb{P}(A)$  as an expectation as follows.

Define the indicator random variable:  $I_A = \begin{cases} 1 & \text{if event } A \text{ occurs,} \\ 0 & \text{otherwise.} \end{cases}$

Then  $\mathbb{E}(I_A) = \mathbb{P}(I_A = 1) = \mathbb{P}(A)$ .

We can refine this expression further, using the idea of conditional expectation. Let  $Y$  be any random variable. Then

continuous  
discrete  
etc

$$\mathbb{P}(A) = \underbrace{\mathbb{E}(I_A)}_{\text{from above}} = \mathbb{E}_Y \left\{ \mathbb{E}(I_A | Y) \right\}$$

↑  
by the Law of Total  
Expectation, true for ANY Y

But  $E(I_A | Y) = \sum_{r=0}^1 r P(I_A = r | Y)$

 $= 0 * P(I_A = 0 | Y) + 1 * P(I_A = 1 | Y)$ 
 $= P(I_A = 1 | Y)$ 

*random variable,  
a function of  
Y (as expected)*

 $= P(A | Y) \text{ by defn of } I_A.$

Thus

$P(A) = E_Y \{ P(A | Y) \}.$

(compare LoTE:  $E(X) = E_Y \{ E(X | Y) \}$ .)

This means that for any random variable X (discrete or continuous), and for any set of values  $S$  (a discrete set or a continuous set), we can write:

and let  $A = \text{any event}$ ,

- for any **discrete** random variable  $Y$ ,

$P(X \in S) = \sum_y P(X \in S | Y=y) P(Y=y)$

or  $P(A) = \sum_y P(A | Y=y) P(Y=y).$

- for any **continuous** random variable  $Y$ ,

$P(X \in S) = \int_y \underbrace{P(X \in S | Y=y)}_{\text{PDF}} f_Y(y) dy$

or  $P(A) = \int_y P(A | Y=y) f_Y(y) dy.$

This IS the Partition Thm:  
we've now learnt it is a special case of the Law of Total Expectation  
now  $f_Y(y)$  is the PDF of  $Y$ .

Interesting because it wasn't clear before that we are ALLOWED to mix P's and pdfs.

Example of probability as a conditional expectation: winning a lottery



Suppose that a million people have bought tickets for the weekly lottery draw. Each person has a probability of one-in-a-million of selecting the winning numbers. If more than one person selects the winning numbers, the winner will be chosen at random from all those with matching numbers.

Winning #s = crime scene DNA sample

one person known to match = criminal

$Y = \# \text{ other unlucky people in AkL who also match by chance}$

Note: Whole analysis takes place in the sample space

$\Omega = \{\text{all draws where I have a matching ticket}\}$

↑  
me: match

○ ○ ○ ○ ○ ○ ----- ○ 1 million

$\Omega = \{\text{everything that can happen to OTHER people}\}$   
 $\Omega = \{\text{these OTHER 1 million people}\}$

You watch the lottery draw on TV and your numbers match the winners!! You had a one-in-a-million chance, and there were a million players, so it must be YOU, right?

Not so fast. Before you rush to claim your prize, let's calculate the probability that you really will win. You definitely win if you are the only person with matching numbers, but you can also win if there are multiple matching tickets and yours is the one selected at random from the matches.

Define  $Y$  to be the number of OTHER matching tickets out of the OTHER 1 million tickets sold. (If you are lucky,  $Y = 0$  so you have definitely won.)

If there are 1 million tickets and each ticket has a one-in-a-million chance of having the winning numbers, then

$Y \sim \text{Poisson}(1)$  approximately.

The relationship  $Y \sim \text{Poisson}(1)$  arises because of the Poisson approximation to the Binomial distribution. In fact,  $Y \sim \text{Binomial}(1\text{million}, \frac{1}{1\text{million}})$

(a) What is the probability function of  $Y$ ,  $f_Y(y)$ ?

$$f_Y(y) = P(Y=y) = \frac{1^y e^{-1}}{y!} = \frac{1}{e * y!} \quad \text{for } y=0, 1, 2, \dots$$

for  $Y \sim \text{Poisson}(1)$

(b) What is the probability that yours is the only matching ticket?

$$P(\text{mine is only match}) = P(Y=0) = \frac{1}{e * 0!} = \frac{1}{e} = 0.368.$$

(c) The prize is chosen at random from all those who have matching tickets.

What is the probability that you win if there are  $Y = y$  OTHER matching tickets?

Let  $W$  be the event that I win.

$$Y=0 \quad P(W|Y=0) = 1$$

$$Y=1 \quad P(W|Y=1) = \frac{1}{2}$$

$$Y=2 \quad P(W|Y=2) = \frac{1}{3}$$

:

$$P(W | Y=y) = \frac{1}{y+1} \quad \begin{matrix} \text{True for} \\ y \text{ others} \rightarrow & \text{me} \uparrow \\ & y+1 \end{matrix} \quad y=0, 1, 2, \dots$$

We don't know what  $Y$  is!

- (d) Overall, what is the probability that you win, given that you have a matching ticket?

$$\begin{aligned}
 P(W) &= \mathbb{E}_Y \{ P(W|Y) \} \\
 &= \sum_{y=0}^{\infty} P(W|Y=y) P(Y=y) \\
 &= \sum_{y=0}^{\infty} \left( \frac{1}{y+1} \right) \left( \frac{1}{e * y!} \right) \\
 &= \frac{1}{e} \sum_{y=0}^{\infty} \frac{1}{(y+1)y!} \\
 &= \frac{1}{e} \sum_{y=0}^{\infty} \frac{1}{(y+1)!} \quad e = \frac{1}{0!} + \frac{1}{1!} + \frac{1}{2!} + \frac{1}{3!} + \dots \\
 &= \frac{1}{e} \left\{ \sum_{y=0}^{\infty} \frac{1}{y!} - \frac{1}{0!} \right\} \\
 &= \frac{1}{e} * (e - 1)
 \end{aligned}$$

using exponential power series for  $e$ .

$$P(W) = 1 - \frac{1}{e}$$

$$= 0.632. \quad \checkmark$$

Disappointing?

See A2 Q4 also

### 3.7 Special process: a model for gene spread

Suppose that a particular gene comes in two variants (alleles): A and B. We might be interested in the case where one of the alleles, say A, is harmful — for example it causes a disease. All animals in the population must have either allele A or allele B. We want to know how long it will take before all animals have the same allele, and whether this allele will be the harmful allele A or the safe allele B. This simple model assumes asexual reproduction. It is very similar to the famous Wright-Fisher model, which is a fundamental model of population genetics.

here we use  $N$  animals  
Wright-Fisher uses  $2N$  genes.

**Assumptions:**

1. The population stays at constant size  $N$  for all generations.
2. At the end of each generation, the  $N$  animals create  $N$  offspring and then they immediately die.  
*N parents,  $\frac{x}{N}$  have gene A,*
3. If there are  $x$  parents with allele A, and  $N - x$  with allele B, then each offspring gets allele A with probability  $x/N$  and allele B with  $1 - x/N$ .  
*then  $P(\text{child has gene A}) = \frac{x}{N}$ .*
4. All offspring are independent.

$| X_t = x$

**Stochastic process:**

The state of the process at time  $t$  is  
 $X_t = \# \text{ of animals with allele A at generation } t.$

Each  $X_t$  could be  $0, 1, 2, \dots, N$ . The state space is  $\{0, 1, 2, \dots, N\}$ .

$N = 12$

$\Rightarrow$  next generation, each child gets  
 $\begin{cases} A & \text{w.p. } 4/12 \\ B & \text{w.p. } 8/12 \end{cases}$

**Distribution of  $[X_{t+1} | X_t]$**

Suppose that  $X_t = x$ , so  $x$  of the animals at generation  $t$  have allele A.

Each of the  $N$  offspring will get A with probability  $\frac{x}{N}$  and B with probability  $1 - \frac{x}{N}$ .

Thus the number of offspring at time  $t+1$  with allele A is:  $X_{t+1} \sim \text{Binomial}(N, \frac{x}{N})$ .

We write this as follows:

$$(X_{t+1} | X_t = x) \sim \text{Binomial}(N, \frac{x}{N}).$$

If

$$[X_{t+1} | X_t = x] \sim \text{Binomial}\left(N, \frac{x}{N}\right),$$

then

$$\Pr(X_{t+1} = y | X_t = x) = \binom{N}{y} \left(\frac{x}{N}\right)^y \left(1 - \frac{x}{N}\right)^{N-y}$$

... to  $y$       arrow going from  $x$  ....      prob

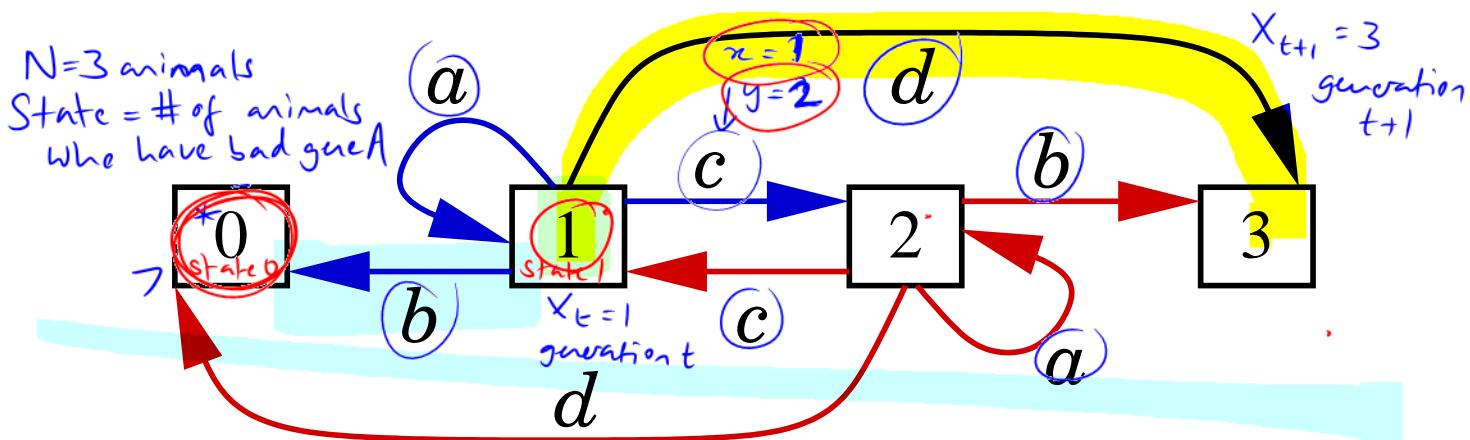
Example with  $N = 3$

This process becomes complicated to do by hand when  $N$  is large. We can use small  $N$  to see how to use first-step analysis to answer our questions.

Google Sirocco the kakapo.

Transition diagram:

**Exercise:** find the missing probabilities  $a, b, c$ , and  $d$  when  $N = 3$ . Express them all as fractions over the same denominator.



Probability the harmful allele A dies out

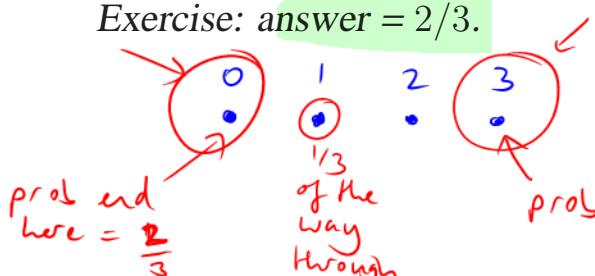
Suppose the process starts at generation 0. One of the three animals has the harmful allele A. Define a suitable notation, and find the probability that the harmful allele A eventually dies out.

use FSA.

Exercise: answer =  $2/3$ .

Ass2 Q4 (d) - (h)

proves this without  
even needing FSA!



prob of ending up here =  $1/3$

## Expected number of generations to fixation

Suppose again that the process starts at generation 0, and one of the three animals has the harmful allele A. Eventually all animals will have the same allele, whether it is allele A or B. When this happens, the population is said to have reached *fixation*: it is fixed for a single allele and no further changes are possible.

Define a suitable notation, and find the expected number of generations to fixation.

*Exercise: answer = 3 generations on average.*

Things get more interesting for large  $N$ . When  $N = 100$ , and  $x = 10$  animals have the harmful allele at generation 0, there is a 90% chance that the harmful allele will die out and a 10% chance that the harmful allele will take over the whole population. The expected number of generations taken to reach fixation is 63.5. If the process starts with just  $x = 1$  animal with the harmful allele, there is a 99% chance the harmful allele will die out, but the expected number of generations to fixation is 10.5. Despite the allele being rare, the *average* number of generations for it to either die out or saturate the population is quite large.

**Note:** The model above is also an example of a process called the **Voter Process**. The  $N$  individuals correspond to  $N$  people who each support one of two political candidates, A or B. Every day they make a new decision about whom to support, based on the amount of current support for each candidate. Fixation in the genetic model corresponds to consensus in the Voter Process.