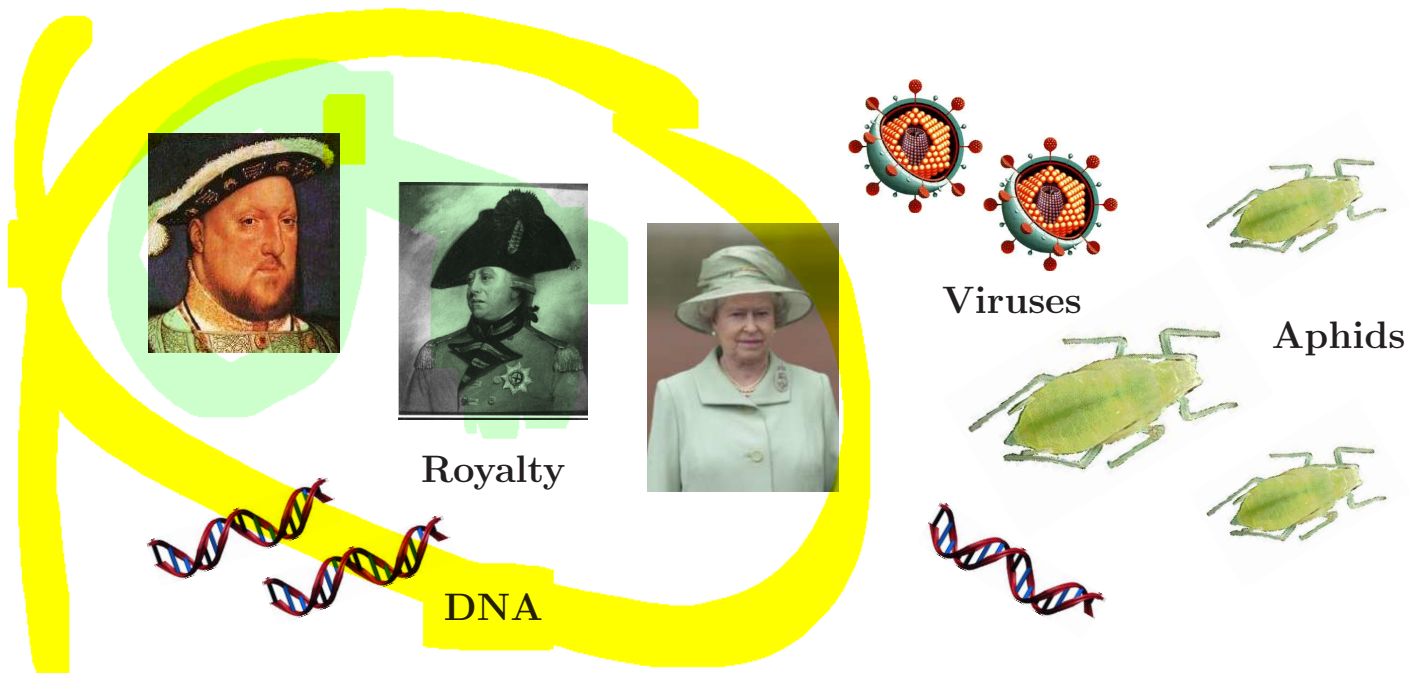


Chapter 6: Branching Processes: The Theory of Reproduction



Although the early development of Probability Theory was motivated by problems in gambling, probabilists soon realised that, if they were to continue as a breed, they must also study

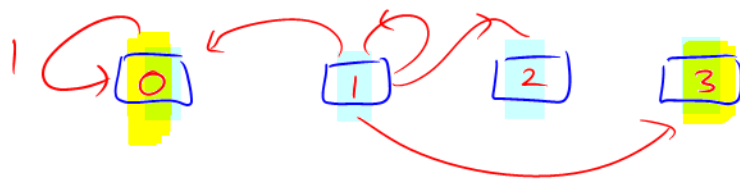
Galton



Reproduction is a complicated business, but considerable insights into population growth can be gained from simplified models. The **Branching Process** is a simple but elegant model of population growth. It is also called the **Galton-Watson Process**, because some of the early theoretical results about the process derive from a correspondence between Sir Francis Galton and the Reverend Henry William Watson in 1873. Francis Galton was a cousin of Charles Darwin. In later life, he developed some less elegant ideas about reproduction — namely eugenics, or selective breeding of humans. Luckily he is better remembered for branching processes.

Rev
Watson





On a transition diagram:
 numbers in the boxes (states)
 ----- are the possible
 values
 of Z_n
 118
 for all n .

Definition: A **branching process** is defined as follows.

- Single individual at time $n = 0$.
- Every individual lives exactly one unit of time, then produces Y offspring, and dies. *Parents and offspring do not coexist.*
- The number of offspring, Y , takes values $0, 1, 2, \dots$, and the probability of producing k offspring is $P(Y=k) = p_k$.
- All individuals reproduce independently. Individuals $1, 2, \dots, n$ have family sizes Y_1, Y_2, \dots, Y_n , where each Y_i has the same distribution as Y .

- • Let Z_n be the number of individuals *{born at time n , for $n = 0, 1, 2, \dots$ Interpret Z_n as the size of generation n .}* *{alive}*
- Then the branching process is $\{Z_0, Z_1, Z_2, Z_3, \dots\} = \{Z_n : n \in \mathbb{N}\}$

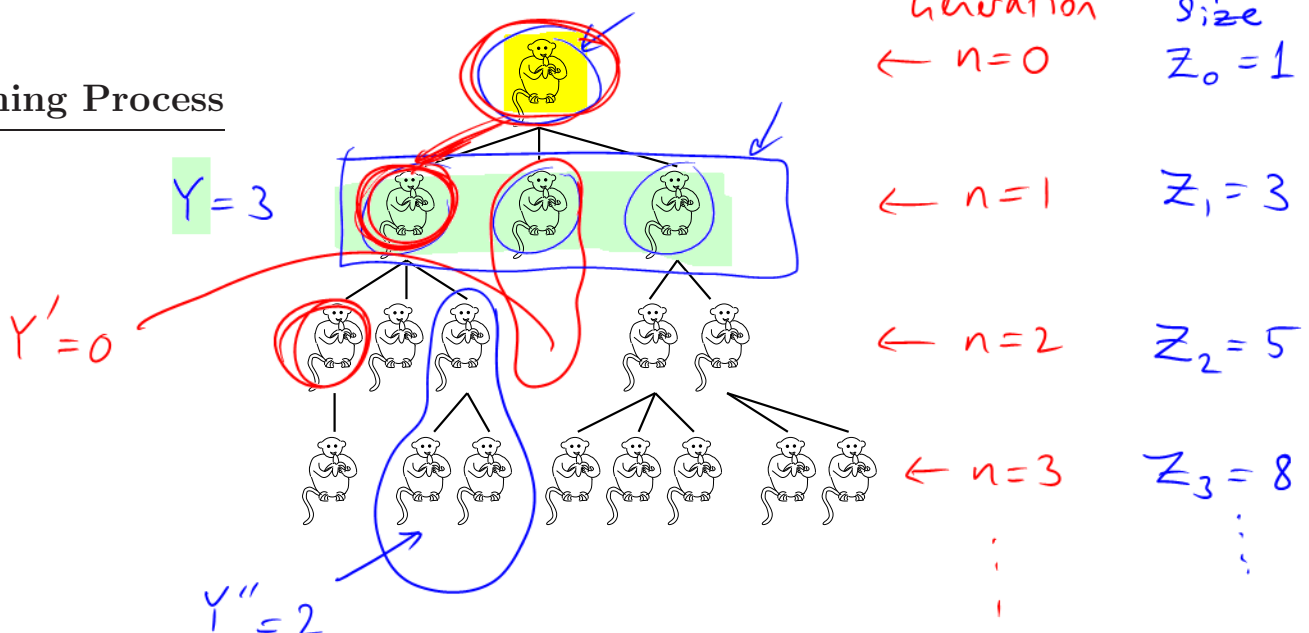
Definition: The *observed* **state** of the branching process at time n is Z_n , where each Z_n can take values $0, 1, 2, \dots$. Note that $Z_0 = 1$ always.
 Z_n is the size of the population at time n .

Note: When we want to say that two random variables X and Y have the same distribution, we write: $X \sim Y$.

For example: $Y_1 \sim Y_2 \sim \dots \sim Y_n \sim Y$.

Note: The definition of the branching process is easily generalized to start with more than one individual at time $n = 0$.

Branching Process



6.2 Questions about the Branching Process

When we have a situation that can be modelled by a branching process, there are several questions we might want to answer.

eg. flu virus

If the branching process is just beginning, what will happen in the future?

1. What can we find out about the distribution of Z_n (the population size at generation n)?

- can we find the mean and variance of Z_n ?

— yes, using the probability generating function of family size, Y ;

fundamental tool: $G(s) = \mathbb{E}(s^Y)$

- can we find the whole distribution of Z_n ?

— for special cases of the family size distribution Y , we can find the PGF of Z_n explicitly;

$Y \sim \text{Geometric}$ is the only non-trivial case.

- can we find the probability that the population has become extinct by generation n , $P(Z_n = 0)$?

— for special cases where we can find the PGF of Z_n (as above).

$Y \sim \text{Geometric}$.

*specific n
↓
makes it hard*

2. What can we find out about eventual extinction?

- can we find the probability of eventual extinction,

— yes, always: using the PGF of Y .

$P(\lim_{n \rightarrow \infty} Z_n = 0)$?

- can we find general conditions for eventual extinction?

— yes: we can find conditions that guarantee that extinction will occur with probability 1.

- if eventual extinction is definite, can we find the distribution of the time to extinction?

— for special cases where we can find the PGF of Z_n (as above).

Example: Modelling cancerous growths. Will a colony of cancerous cells become extinct before it is sufficiently large to overgrow the surrounding tissue?

Most Recent Common Ancestor ?

5000 - 2000 years ago



If the branching process is already in progress, what happened in the past?

1. How long has the process been running?

- how many generations do we have to go back to get to the single common ancestor?

2. What has been the distribution of family size over the generations?

3. What is the total number of individuals (over all generations) up to the present day?
#humans ever born ~ 106 billion starting 50000 years ago.

Example: It is believed that all humans are descended from a single female ancestor, who lived in Africa. How long ago?

Estimated at ~200,000 years.

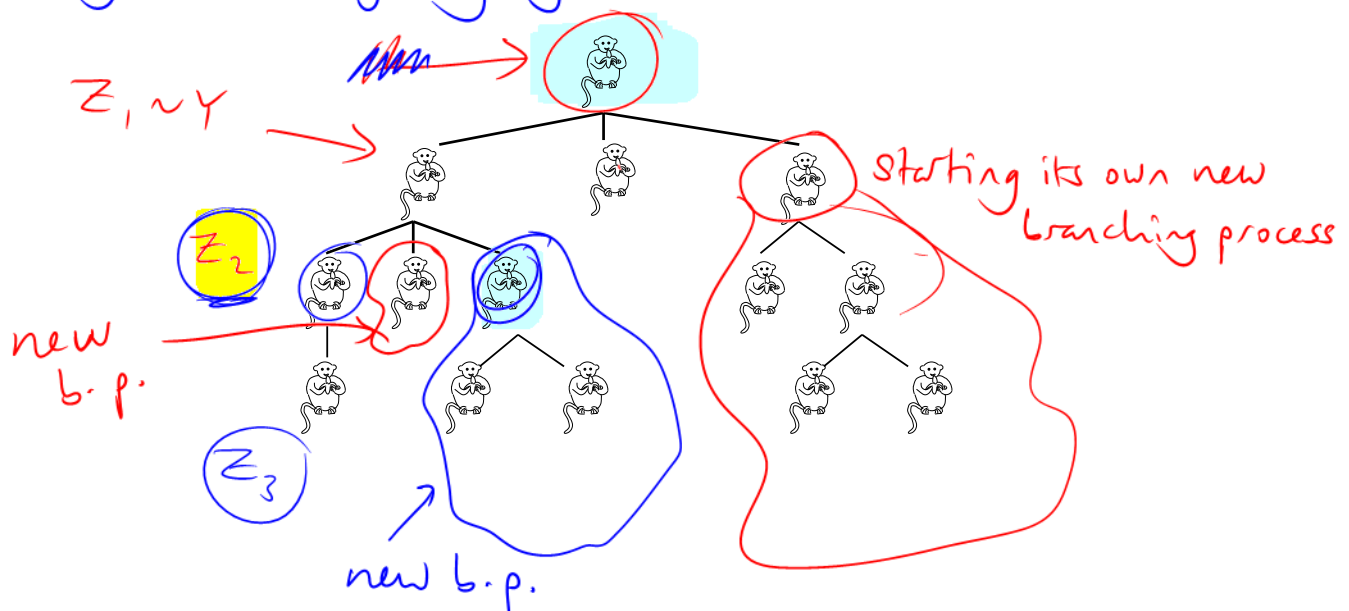
What has been the mean family size over that period?

- Very close to 1 female offspring per female adult;
estimate = 1.002.



6.3 Analysing the Branching Process

Key Observation: every individual in every generation starts a new, independent branching process, as if the whole process were starting at the beginning again.



state, size @ gen n
 Z_n as a randomly stopped sum

Most of the interesting properties of the branching process centre on the distribution of Z_n (the population size at time n). Using the Key Observation from overleaf, we can find an expression for the probability generating function of Z_n .

Consider the following.

- The population size at time $n-1$ is given by Z_{n-1} .
(think of as Z_{n-1} parents)
- Label the individuals at time $n-1$ as $1, 2, 3, \dots, Z_{n-1}$.
- Each parent $1, 2, \dots, Z_{n-1}$ starts a new branching process. Let $Y_1, Y_2, \dots, Y_{Z_{n-1}}$ be the random family sizes of the parents $1, 2, \dots, Z_{n-1}$.
randomly stopped because we don't yet know how many parents at time $n-1$.
- The #offspring at time n , Z_n , is equal to the sum of offspring of the Z_{n-1} parents.

That is:

$$Z_n = \sum_{i=1}^{Z_{n-1}} Y_i$$

randomly stopped sum!

So Z_n is a randomly stopped sum:

a sum of Y_1, Y_2, \dots randomly stopped by the random variable Z_{n-1} . [which is unknown at the point we start from]

Note: 1. Each $Y_i \sim Y$.

2. $Y_1, Y_2, \dots, Y_{Z_{n-1}}$ are independent of each other and of Z_{n-1} .

Thm 4.6: $T_N = X_1 + \dots + X_N$ $N \sim \text{random}$.

$$\Rightarrow G_{T_N}(s) = G_N(G_X(s))$$

Here: $Z_n = Y_1 + \dots + Y_{Z_{n-1}}$ $Z_{n-1} \sim \text{random}$
 $\Rightarrow G_{Z_n}(s) = G_{Z_{n-1}}(G_Y(s))$.

Probability Generating Function of Z_n

Let $G_Y(s) = \mathbb{E}(s^Y)$ be the probability generating function of Y .

(Recall that Y is the number of Young of an individual: the family size.

children @ time n

Now Z_n is a randomly stopped sum: it is the sum of Y_1, Y_2, \dots , stopped by the random variable Z_{n-1} . So we can use Theorem 4.6 (Chapter 4) to express the PGF of Z_n directly in terms of the PGFs of Y and Z_{n-1} .

parents @ time $n-1$.

By Theorem 4.6, if $Z_n = Y_1 + Y_2 + \dots + Y_{Z_{n-1}}$, and Z_{n-1} is itself random, then the PGF of Z_n is given by:

$$G_{Z_n}(s) = G_{Z_{n-1}}(G_Y(s)) \quad (*)$$

where $G_{Z_{n-1}}$ is the PGF of the random variable Z_{n-1} .

For ease of notation, we can write:

$$G_{Z_n}(s) = G_n(s), \quad G_{Z_{n-1}}(s) = G_{n-1}(s), \quad \text{and so on.}$$

Note that $Z_1 \sim Y$ (the # individuals born at time $n=1$) so we can also write:

$$G_Y(s) = G_1(s) = G(s) \text{ for simplicity.}$$

Thus, from $(*)$,

$$G_n(s) = G_{n-1}(G(s))$$

Branching process recursion formula

our fundamental repeating unit
 \rightarrow make notation as easy as possible.

Note:

1. $G_n(s) = \mathbb{E}(s^{Z_n})$, the PGF of the population size at time n , Z_n .
2. $G_{n-1}(s) = \mathbb{E}(s^{Z_{n-1}})$, the PGF of popn size at time $n-1$, Z_{n-1} .
3. $G(s) = \mathbb{E}(s^Y) = \mathbb{E}(s^{Z_1}) = G_1(s)$, the PGF of family size, Y .

We are trying to find the PGF of Z_n , the population size at time n .

So far, we have:
$$G_n(s) = G_{n-1}(G(s)). \quad (*)$$

*r & s are
just variable
names*

But by the same argument,

$$G_{n-1}(r) = G_{n-2}(G(r)) \quad (**)$$

(using r instead of s to avoid confusion in the next line).

Substituting in (*),

$$\begin{aligned} G_n(s) &= G_{n-1}(G(s)) \quad (*) \\ &= G_{n-1}(r) \quad \text{put } r = G(s) \\ &= G_{n-2}(G(r)) \quad \text{by } (**) \\ &= G_{n-2}(G(G(s))) \quad \text{replacing } r = G(s) \end{aligned}$$

By the same reasoning, we will obtain:

$$G_n(s) = G_{n-3}(G(G(G(s))))$$

n-3 *3 times*

and so on, until we finally get:

$$G_n(s) = G_{n-(n-1)}(G(G(\dots G(s) \dots)))$$

n-(n-1) *n-1 times*

but
 $G_1 = G$

$$\begin{aligned} &\xrightarrow{\quad} G_1(G(G(\dots G(s) \dots))) \\ \therefore G_n(s) &= G(G(G(\dots G(s) \dots))) \end{aligned}$$

n times.

We have therefore proved the following Theorem.

Theorem 6.3: Let $G(s) = \mathbb{E}(s^Y) = \sum_{y=0}^{\infty} p_y s^y$ be the PGF of the family size distribution, Y . Let $Z_0 = 1$ (start from a single individual at time 0), and let Z_n be the population size at time n ($n = 0, 1, 2, \dots$). Let $G_n(s)$ be the PGF of the random variable Z_n . Then

$$G_n(s) = G\left(\underbrace{G\left(G\left(\dots G(s)\dots\right)\right)}_{n \text{ times}}\right). \quad \square$$

Note: $G_n(s) = \underbrace{G\left(G\left(G\left(\dots G(s)\dots\right)\right)\right)}_{n \text{ times}}$ is called the n -fold iterate of G .

We have therefore found an expression for the PGF of the population size at generation n , although there is no guarantee that it is possible to write it down or manipulate it very easily for large n . For example, if Y has a $\text{Poisson}(\lambda)$ distribution, then $G(s) = e^{\lambda(s-1)}$, and already by generation $n = 3$ we have the following fearsome expression for $G_3(s)$:

$$G_3(s) = e^{\lambda\left(e^{\lambda\left(e^{\lambda(s-1)}-1\right)}-1\right)}. \quad (\text{Or something like that!})$$

However, in some circumstances we can find quite reasonable closed-form expressions for $G_n(s)$, notably when Y has a Geometric distribution. In addition, for any distribution of Y we can use the expression $G_n(s) = G_{n-1}(G(s))$ to derive properties such as the mean and variance of Z_n , and the probability of eventual extinction ($\mathbb{P}(Z_n = 0)$ for some n).

6.4 What does the distribution of Z_n look like?

Before deriving the mean and the variance of Z_n , it is helpful to get some intuitive idea of how the branching process behaves. For example, it seems reasonable to calculate the mean, $\mathbb{E}(Z_n)$, to find out what we expect the population size to be in n generations time, but why are we interested in $\text{Var}(Z_n)$?

The answer is that Z_n usually has a “boom-or-bust” distribution: either the population will take off (boom), and the population size grows quickly, or the population will fail altogether (bust). In fact, if the population fails, it is likely to do so very quickly, within the first few generations. This explains why we are

interested in $\text{Var}(Z_n)$. A huge variance will alert us to the fact that the process does not cluster closely around its mean values. In fact, the mean might be almost useless as a measure of what to expect from the process.

$Y=0$ possible
Simulation 1: $Y \sim \text{Geometric}(p = 0.3)$

mean of Y : $E(Y) = \frac{2}{p} = \frac{0.7}{0.3} = 2.3$

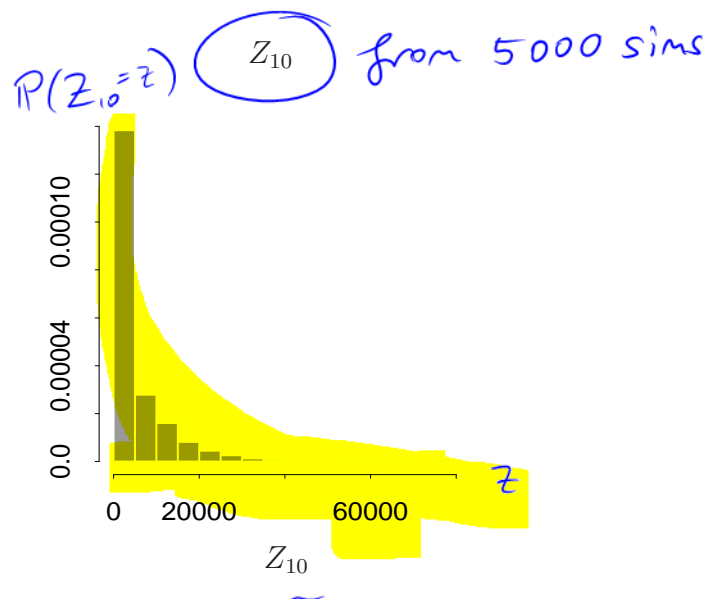
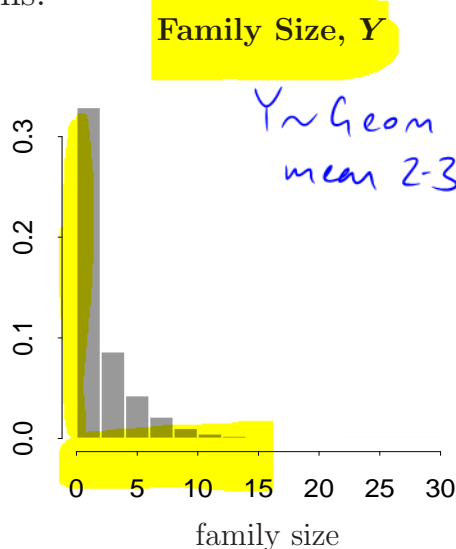
The following table shows the results from 10 simulations of a branching process, where the family size distribution is $Y \sim \text{Geometric}(p = 0.3)$.

10 sims

Simulation	Z_0	Z_1	Z_2	Z_3	Z_4	Z_5	Z_6	Z_7	Z_8	Z_9	Z_{10}	
1	1	0	0	0	0	0	0	0	0	0	0	<i>Extinct.</i>
2	1	1	0	0	0	0	0	0	0	0	0	<i>Extinct at gen 2.</i>
3	1	4	19	42	81	181	433	964	2276	5383	12428	
4	1	3	3	5	3	15	29	86	207	435	952	
5	1	0	0	0	0	0	0	0	0	0	0	
6	1	1	0	0	0	0	0	0	0	0	0	
7	1	2	8	26	68	162	360	845	2039	4746	10941	
8	1	1	0	0	0	0	0	0	0	0	0	
9	1	1	0	0	0	0	0	0	0	0	0	
10	1	1	4	13	18	39	104	294	690	1566	3534	

Often, the population is extinct by generation 10. However, when it is not extinct, it can take enormous values (12428, 10941, ...).

The same simulation was repeated 5000 times to find the empirical distribution of the population size at generation 10 (Z_{10}). The figures below show the distribution of family size, Y , and the distribution of Z_{10} from the 5000 simulations.



In this example, the family size is rather variable, but the variability in Z_{10} is enormous (note the range on the histogram from 0 to 60,000). Some statistics are:

Proportion of samples extinct by generation 10: 0.436

Summary of Z_n :

Min	1st Qu	Median	Mean	3rd Qu	Max
0	0	1003	4617	6656	82486

Mean of Z_n :

4617.2

Variance of Z_n :

53937785.7

theory $\Rightarrow 4784$

Huge variance \Rightarrow mean is virtually uninformative.

So the empirical variance is $\text{Var}(Z_{10}) = 5.39 \times 10^7$. This perhaps contains more useful information than the mean value of 4617. The distribution of Z_n has 43.6% of zeros, but (when it is non-zero) takes values up to 82,486. Is it really useful to summarize such a distribution by the single mean value 4617?

For interest, out of the 5000 simulations, there were only 35 (0.7%) that had a value for Z_{10} greater than 0 but less than 100. This emphasizes the “boom-or-bust” nature of the distribution of Z_n .

Simulation 2: $Y \sim \text{Geometric}(p = 0.5)$

$$\leftarrow EY = \frac{q}{p} = \frac{0.5}{0.5} = 1.$$

We repeat the simulation above with a different value for p in the Geometric family size distribution: this time, $p = 0.5$. The family size distribution is therefore

Simulation	Z_0	Z_1	Z_2	Z_3	Z_4	Z_5	Z_6	Z_7	Z_8	Z_9	Z_{10}
1	1	0	0	0	0	0	0	0	0	0	0
2	1	0	0	0	0	0	0	0	0	0	0
3	1	0	0	0	0	0	0	0	0	0	0
4	1	0	0	0	0	0	0	0	0	0	0
5	1	1	0	0	0	0	0	0	0	0	0
6	1	7	9	17	15	20	19	8	7	13	35
7	1	2	5	2	5	8	8	3	3	0	0
8	1	2	0	0	0	0	0	0	0	0	0
9	1	0	0	0	0	0	0	0	0	0	0
10	1	0	0	0	0	0	0	0	0	0	0

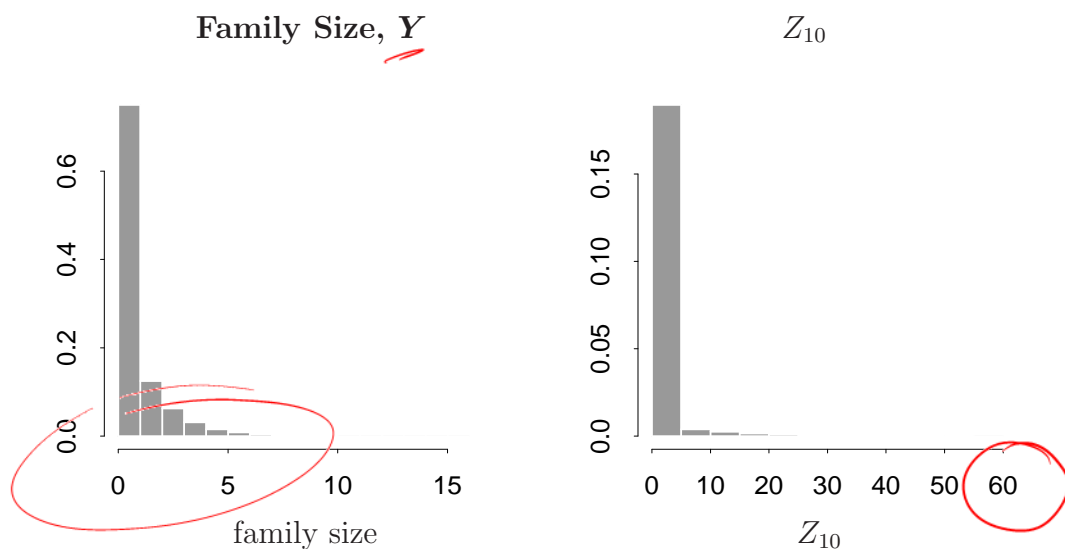
10
sim

\rightarrow

This time, almost all the populations become extinct. We will see later that this value of p (just) guarantees eventual extinction with probability 1.

but $E(\text{time taken}) = \infty$.

The family size distribution, $Y \sim \text{Geometric}(p = 0.5)$, and the results for Z_{10} from 5000 simulations, are shown below. Family sizes are often zero, but families of size 2 and 3 are not uncommon. It seems that this is not enough to save the process from extinction. This time, the maximum population size observed for Z_{10} from 5000 simulations was only 56, and the mean and variance of Z_{10} are much smaller than before.



Proportion of samples extinct by generation 10: 0.9108

Summary of Z_n :

Min	1st Qu	Median	Mean	3rd Qu	Max
0	0	0	0.965	0	56

Mean of Z_n : 0.965

Variance of Z_n : 19.497

What happens for larger values of p ?

It was mentioned above that $Y \sim \text{Geometric}(p = 0.5)$ just guarantees eventual extinction with probability 1. For $p > 0.5$, extinction is also guaranteed, and tends to happen quickly. For example, when $p = 0.55$, over 97% of simulated populations are already extinct by generation 10.

6.5 Mean and variance of Z_n

The previous section has given us a good idea of the significance and interpretation of $\mathbb{E}(Z_n)$ and $\text{Var}(Z_n)$. We now proceed to calculate them. Both $\mathbb{E}(Z_n)$ and $\text{Var}(Z_n)$ can be expressed in terms of *the mean and variance of the family size distribution, Y .*

Thus, let $\mathbb{E}Y = \mu$ and let $\text{Var}(Y) = \sigma^2$. These are the mean & variance of the #offspring of a single individual.

Theorem 6.5: Let $\{Z_0, Z_1, Z_2, \dots\}$ be a branching process with $Z_0 = 1$ (start with a single individual). Let Y denote the family size distribution, and suppose that $\mathbb{E}(Y) = \mu$. Then

$$\mathbb{E}(Z_n) = \mu^n.$$

Proof:

By p. 121, $Z_n = Y_1 + Y_2 + \dots + Y_{Z_{n-1}}$

$$Z_n = \sum_{i=1}^{Z_{n-1}} Y_i.$$

Thus, from Section 3.4 (p. 62):

$$\begin{aligned} \mathbb{E}(Z_n) &= \mathbb{E}(Y_i) * \mathbb{E}(Z_{n-1}) \\ &= \mu \mathbb{E}(Z_{n-1}) \\ &= \mu \{ \mu \mathbb{E}(Z_{n-2}) \} \\ &= \mu^2 \mathbb{E}(Z_{n-2}) \\ &\vdots \\ &= \mu^{n-1} \mathbb{E}(Z_{n-(n-1)}) \\ &= \mu^{n-1} \mathbb{E}(Z_1) \end{aligned}$$

$$\therefore \mathbb{E}(Z_n) = \mu^n$$

□

e.g. $Y=2$ w.p. 1
 $Z_0 = 1 = 2^0$
 $Z_1 = 2 = 2^1$
 $Z_2 = 2^2$
 $Z_3 = 2^3$ etc.

is a randomly stopped sum:

p. 62:
 $T_N = X_1 + \dots + X_N$
 with $\mathbb{E}(X_i) = \mu$
 $\Rightarrow \mathbb{E}(T_N) = \mu \mathbb{E}N$ ←
 Here, $T_N = Z_n$
 and $X_i = Y_i$ and $N = Z_{n-1}$

but $\mathbb{E}Z_1 = \mathbb{E}Y = \mu$

Examples: Consider the simulations of Section 6.4.

1. Family size $Y \sim \text{Geometric}(p = 0.3)$. So $\mu = \mathbb{E}Y = \frac{q}{p} = \frac{0.7}{0.3} = 2.33$.

Expected population size by generation $n = 10$ is:

$$\mathbb{E}(Z_{10}) = \mu^{10} = (2.33)^{10} = 4784.$$

The theoretical value, 4784, compares well with the sample mean from 5000 simulations, 4617 (page 126).

2. Family size $Y \sim \text{Geometric}(p = 0.5)$. So $\mu = \mathbb{E}Y = \frac{q}{p} = \frac{0.5}{0.5} = 1$,

$$\text{so } \mathbb{E}(Z_{10}) = \mu^{10} = 1^{10} = 1.$$

Compares well with the sample mean of 0.965 (page 127).

Variance of Z_n

Arguments very similar to Ass 4 Q3.

Theorem 6.5: Let $\{Z_0, Z_1, Z_2, \dots\}$ be a branching process with $Z_0 = 1$ (start with a single individual). Let Y denote the family size distribution, and suppose that $\mathbb{E}(Y) = \mu$ and $\text{Var}(Y) = \sigma^2$. Then

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

Proof:

Write $V_n = \text{Var}(Z_n)$. The proof works by finding a recursive formula for V_n .

p.62: $\text{Var}(T_N) = \sigma^2 \mathbb{E}N + \mu^2 \text{Var}(N)$

Here, $N = Z_{n-1}$

so $\mathbb{E}N = \mathbb{E}(Z_{n-1}) = \mu^{n-1}$

and $\text{Var}(N) = \text{Var}(Z_{n-1}) = V_{n-1}$

Define $V_n = \text{Var}(Z_n)$

Using the Law of Total Variance for randomly stopped sums from Section 3.4 (page 62),

$$Z_n = \sum_{i=1}^{Z_{n-1}} Y_i$$

$$\Rightarrow \text{Var}(Z_n) = \{\mathbb{E}(Y_i)\}^2 \times \text{Var}(Z_{n-1}) + \text{Var}(Y_i) \times \mathbb{E}(Z_{n-1})$$

$$\Rightarrow V_n = \mu^2 V_{n-1} + \sigma^2 \mathbb{E}(Z_{n-1})$$

$$\Rightarrow V_n = \mu^2 V_{n-1} + \sigma^2 \mu^{n-1}, \quad \dots$$

using $\mathbb{E}(Z_{n-1}) = \mu^{n-1}$ as above.

Also,

$$V_1 = \text{Var}(Z_1) = \text{Var}(Y) = \sigma^2.$$

Find V_n by repeated substitution:

$$V_1 = \sigma^2$$

$$V_2 = \mu^2 V_1 + \sigma^2 \mu = \mu^2 \sigma^2 + \mu \sigma^2 = \mu \sigma^2 (1 + \mu)$$

$$V_3 = \mu^2 V_2 + \sigma^2 \mu^2 = \mu^2 \sigma^2 (1 + \mu + \mu^2)$$

$$V_4 = \mu^2 V_3 + \sigma^2 \mu^3 = \mu^3 \sigma^2 (1 + \mu + \mu^2 + \mu^3)$$

\vdots etc.

Completing the pattern,

$$V_n = \mu^{n-1} \sigma^2 (1 + \mu + \mu^2 + \dots + \mu^{n-1})$$

$$= \mu^{n-1} \sigma^2 \sum_{r=0}^{n-1} \mu^r$$

$$= \mu^{n-1} \sigma^2 \left(\frac{1 - \mu^n}{1 - \mu} \right).$$

Valid for $\mu \neq 1$.

(sum of first n terms of Geometric series)

prove this sum in Ass 4 Q 3a

When $\mu = 1$:

$$V_n = 1^{n-1} \sigma^2 \underbrace{(1^0 + 1^1 + \dots + 1^{n-1})}_{n \text{ times}} = \sigma^2 n.$$

Hence the result:

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

Examples: Again consider the simulations of Section 6.4.

1. Family size $Y \sim \text{Geometric}(p = 0.3)$. So $\mu = \mathbb{E}(Y) = \frac{q}{p} = \frac{0.7}{0.3} = 2.33$.
 $\sigma^2 = \text{Var}(Y) = \frac{q}{p^2} = \frac{0.7}{(0.3)^2} = 7.78$.

$$\text{So } \text{Var}(Z_{10}) = \sigma^2 \mu^9 \left(\frac{1 - \mu^{10}}{1 - \mu} \right) = 5.72 \times 10^7.$$

Compares well with the sample variance from 5000 simulations, 5.39×10^7 (page 126).

2. Family size $Y \sim \text{Geometric}(p = 0.5)$. So $\mu = \mathbb{E}(Y) = \frac{q}{p} = \frac{0.5}{0.5} = 1$.

$$\sigma^2 = \text{Var}(Y) = \frac{q}{p^2} = \frac{0.5}{(0.5)^2} = 2. \quad \text{Using the formula for } \text{Var}(Z_n) \text{ when } \mu=1, \text{ we get } \text{Var}(Z_{10}) = \sigma^2 n = 2 \times 10 = 20.$$

Compares well with the sample variance of 19.5 (page 127).