# Chapter 8: Branching Processes:

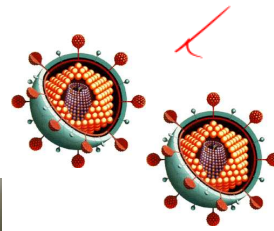## The Theory of Reproduction



Viruses

Aphids

Royalty

DNA

*Galton*    ♂

*Rev Watson*

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



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.

*Journal* Annals of Human Genetics

→ *Eugen.* Annals of Eugenics.

## 8.1 Branching Processes

Consider some sort of **population** consisting of reproducing individuals.

**Examples:** living things (animals, plants, bacteria, royal families);
diseases; computer viruses;
rumours, gossip, lies (one lie always leads to another!)

**Start conditions:** start at time $n=0$, with a single individual.

**Each individual:** lives for 1 unit of time. At time $n=1$, it produces a family of offspring, and immediately dies.

**How many offspring?** Could be $0, 1, 2, \ldots$ This is the family size, $\underline{\text{size}}$, $Y$. [ $Y$ stands for "number of Young". ]

**Each offspring:** lives for 1 unit of time. At time $n=2$, it produces its own family of offspring, and immediately dies.

**and so on...**

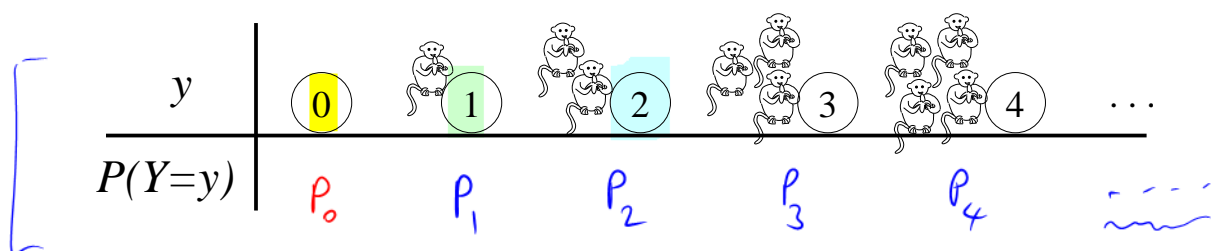### Assumptions

1. All individuals reproduce independently of each other.

2. The family sizes of different individuals are independent, identically distributed r.v.s. We denote the family size by $Y$ = number of Young.

**Family size distribution, $Y$**    $P(Y=k) = p_k$.

| $y$ | 0 | 1 | 2 | 3 | 4 | $\cdots$ |
|---|---|---|---|---|---|---|
| $P(Y=y)$ | $p_0$ | $p_1$ | $p_2$ | $p_3$ | $p_4$ | |

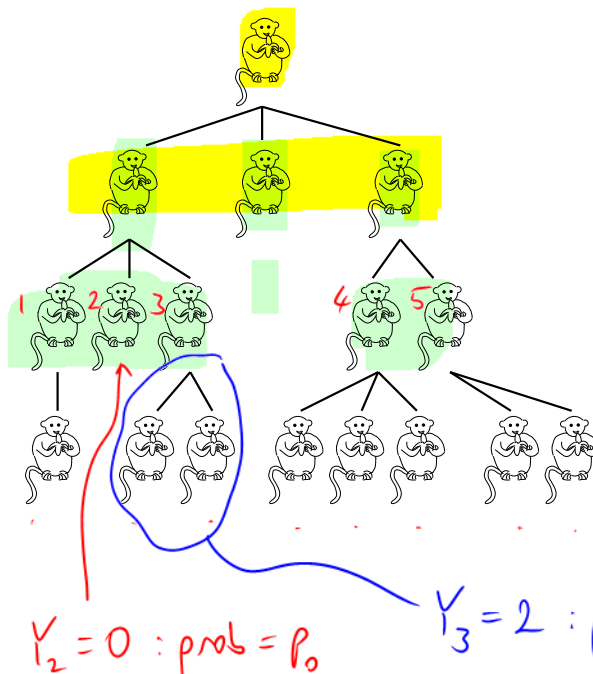*Definition:* A **branching process** is defined as follows.

- Single individual at time $n = 0$.

- Every individual lives exactly one unit of time, then produces $Y$ offspring, and dies.    *Parents & Offspring never coexist.*

- The number of offspring, $Y$, takes values 0, 1, 2, ..., and the probability of producing $k$ offspring is    $P(Y=k) = p_k$ .

- All individuals reproduce independently. Individuals $1, 2, \ldots, n$ have family sizes $Y_1, Y_2, \ldots Y_n$, where *each $Y_i \sim Y$. ← ie. each $Y_i$ has same distn as $Y$.*

- Let $Z_n$ be the *number of individuals born at time $n$, for $n = 0, 1, 2, \ldots$ Interpret $Z_n$ as the size of generation $n$.*

- Then the branching process is $\{Z_0, Z_1, Z_2, Z_3, \ldots\} = \{Z_n : n \in \mathbb{N}\}$.

*Definition:* The **state** of the branching process at time $n$ is $z_n$, *where each $z_n$ can take values 0, 1, 2, ... ie. the State Space $= \{0, 1, 2, \ldots\}$. Note that $z_0 = 1$ always.*

***Note:*** When we want to say that two random variables $X$ and $Y$ have the same distribution, we write: $X \sim Y$.
For example: $Y_i \sim Y$, *where $Y_i$ is the family size of individual $i$.*

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

**Branching Process**



| Generation | Popn Size |
|---|---|
| 0 | $Z_0 = 1$ ← |
| 1 | $Z_1 = 3$ |
| 2 | $Z_2 = 5$ |
| 3 | $Z_3 = 8$ |

$Y_2 = 0$ : prob $= p_0$

$Y_3 = 2$ : prob $= p_2$

## 8.2  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.

*e.g. disease like Swine Flu*

**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, $Y$;*

    - 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;*    *$Y \sim$ Geometric is only non-trivial case.*

    - can we find the probability that the population has become **extinct** by generation $n$, $\mathbb{P}(Z_n = 0)$ ?
      — *for special cases where we can find the PGF of $Z_n$ (as above).*    *$Y \sim$ Geometric*

2.  What can we find out about <u>eventual</u> extinction?

    - can we find the <u>probability</u> of eventual extinction, $\mathbb{P}\left(\lim_{n \to \infty} Z_n = 0\right)$ ?
      — *yes, always, using the PGF of $Y$.*    hitting probability !!

    - can we find general <u>conditions</u> for eventual extinction?
      — *yes, we can find conditions that guarantee extinction will occur with probability 1.*    $\mathbb{E} Y \leq 1$

    - 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 $Z_n$.*

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

# 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?

**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 ago.

What has been the mean family size over that period?

— probably very close to 1 female offspring per female adult, e.g. estimate 1·002.
i.e. for every 500 female parents, we get one extra female offspring !!

## 8.3 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.

# Most recent common ancestor?

~3000 years ago.

Estimated 2000 − 5000



most recent common ancestor

— everyone alive today

Pop size : time 0 $\qquad$ size 1

time 1 $\qquad$ $Z_1 \sim Y$

time 2 $\qquad$ $Z_2$

$\vdots$

time $n-1$ $\qquad$ $Z_{n-1}$ parents for time $n$

time $n$ $\qquad$ $Z_n$ offspring at time $n$

## $Z_n$ 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 $Z_{n-1}$.
  ( = #parents at time $n-1$, that will produce generation $n$.)

- Label the individuals at time $n-1$ (the parents) as
  $1, 2, 3, \cdots, Z_{n-1}$.
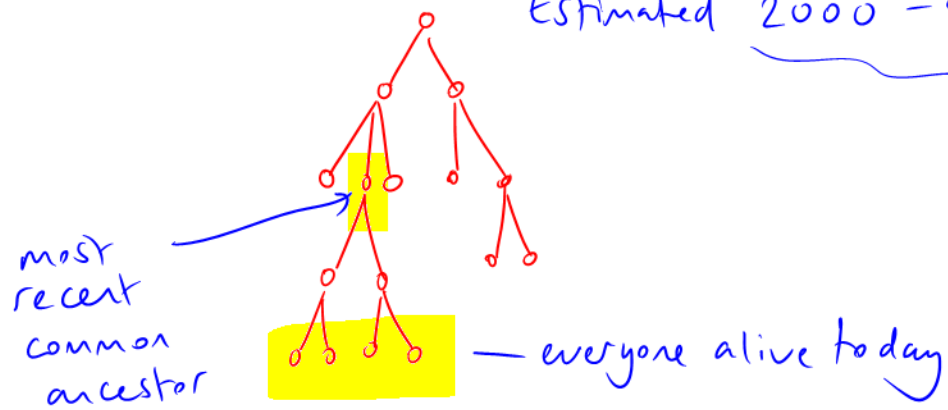
- Let $Y_1, Y_2, \ldots, Y_{Z_{n-1}}$ be the random family sizes of the parents $1, 2, \cdots, Z_{n-1}$.

- The # offspring at time $n$, $Z_n$, is equal to the total number of offspring of the parents $1, 2, \cdots, Z_{n-1}$.

That is, $$Z_n = \sum_{i=1}^{Z_{n-1}} Y_i \quad \leftarrow$$

| time $n-1$ | 1 | 2 | $\cdots$ | $Z_{n-1}$ | parents |
|---|---|---|---|---|---|
| $n$ | $Y_1$ | $Y_2$ | $\cdots$ | $Y_{Z_{n-1}}$ | children |

total #children $= Y_1 + Y_2 + \cdots + Y_{Z_{n-1}}$

So $Z_n$ is a Randomly Stopped Sum: a sum of $Y_1, Y_2, \ldots,$ randomly stopped by the r.v. $Z_{n-1}$.

**Note:** 1. Each $Y_i \sim Y$ : ie. each individual $i = 1, \ldots, Z_{n-1}$ has the same family size distribution.

2. $Y_1, Y_2, \ldots, Y_{Z_{n-1}}$ are independent.

RSS : $Z_n = Y_1 + \cdots + Y_{Z_{n-1}}$

$$G_{Z_n}(s) = \mathbb{E}(s^{Z_n}) = \mathbb{E}_{Z_{n-1}}\left\{ \mathbb{E}(s^{Y_1} \cdots s^{Y_{Z_{n-1}}} | Z_{n-1}) \right\}$$

$$[\text{indep}] \Rightarrow = \mathbb{E}_{Z_{n-1}}\left\{ \mathbb{E}(s^{Y_1}) \cdots \mathbb{E}(s^{Y_{Z_{n-1}}}) \right\}$$

$$= \mathbb{E}_{Z_{n-1}}\left\{ \mathbb{E}(s^Y)^{Z_{n-1}} \right\} = G_{Z_{n-1}}\left( G_Y(s) \right) \quad *$$

THE UNIVERSITY OF AUCKLAND NEW ZEALAND Te Whare Wānanga o Tāmaki Makaurau    168

## 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 # **Young** of one individual : the **family size**.)

**# children**

Now $Z_n$ is a randomly stopped sum: it is the sum of $Y_1, Y_2, \ldots$, stopped by the random variable $Z_{n-1}$. **# parents**. So we can use Theorem 7.6 (Chapter 7) to express the PGF of $Z_n$ directly in terms of the PGFs of $Y$ and $Z_{n-1}$.

By Theorem 7.6, if $Z_n = Y_1 + Y_2 + \ldots + Y_{Z_{n-1}}$, and $Z_{n-1}$ is itself random, then the PGF of $Z_n$ is given by:

$$G_{Z_n}(s) = G_{Z_{n-1}}\left( G_Y(s) \right) \qquad \text{⊛}$$

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

For ease of notation, we can write:

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

Note that   $Z_1 \sim Y$   ( # offspring of a single individual, the **single parent at time n=0**)

so we can also write:

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

Thus,  ⊛  implies ;

$$\boxed{G_n(s) = G_{n-1}\left( G(s) \right).} \qquad \text{Branching Process Recursion Formula}$$

**Note:**

1.  $G_n(s) = \mathbb{E}(s^{Z_n})$, the PGF of population size at time $n$, $Z_n$.

2.  $G_{n-1}(s) = \mathbb{E}(s^{Z_{n-1}})$, " " " " " " " $n-1$, $Z_{n-1}$.

3.  $G(s) = \mathbb{E}(s^Y) = \mathbb{E}(s^{Z_1})$, the PGF of family size, $Y$.

   ( $G(s) = G_1(s)$ by definition.)

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

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

But by the same argument,

$$G_{n-1}(r) = G_{n-2}\Big(G(r)\Big) \qquad (\ast\ast)$$

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

Substituting in $(\star)$,

$$G_n(s) = G_{n-1}\big(G(s)\big)$$
$$= G_{n-1}(r) \qquad \text{where } r = G(s)$$
$$= G_{n-2}\big(G(r)\big) \qquad \text{using } (\ast\ast)$$
$$= G_{n-2}\big(G(G(s))\big) \qquad \text{replacing } r = G(s).$$

By the same reasoning, we will obtain:

$$G_n(s) = G_{n-3}\Big(\underbrace{G(G(G(s)))}_{3\ G's}\Big)$$
$$\underbrace{\phantom{G_{n-3}}}_{n-3}$$

and so on, until we finally get:

$$G_n(s) = G_{n-(n-1)}\Big(\underbrace{G(G(G\ldots G(s)\ldots))}_{n-1\ times}\Big)$$

$$= G_1\Big(\underbrace{G(G\ldots G(s))\ldots)}_{n-1\ times}\Big)$$

but
$G_1 = G$
by definition

$$= \underbrace{G(G(G(\ldots G(s)\ldots)))}_{n\ times.}$$

We have therefore proved the following Theorem.

**Theorem 8.3:** 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 $n$ $(n = 0, 1, 2, \ldots)$. Let $G_n(s)$ be the PGF of the random variable $Z_n$. Then

$$G_n(s) = \underbrace{G\Big(G\Big(G\Big(\ldots G(s) \ldots\Big)\Big)\Big)}_{n \text{ times}}. \qquad \square$$

*Note:* $G_n(s) = \underbrace{G\Big(G\Big(G\Big(\ldots G(s) \ldots\Big)\Big)\Big)}_{n \text{ times}}$ 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 Poisson($\lambda$) distribution, then $G(s) = e^{\lambda(s-1)}$, and already by generation $n = 3$ we have the following fearsome expression for $G_3(s)$:

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

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

## 8.4 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.
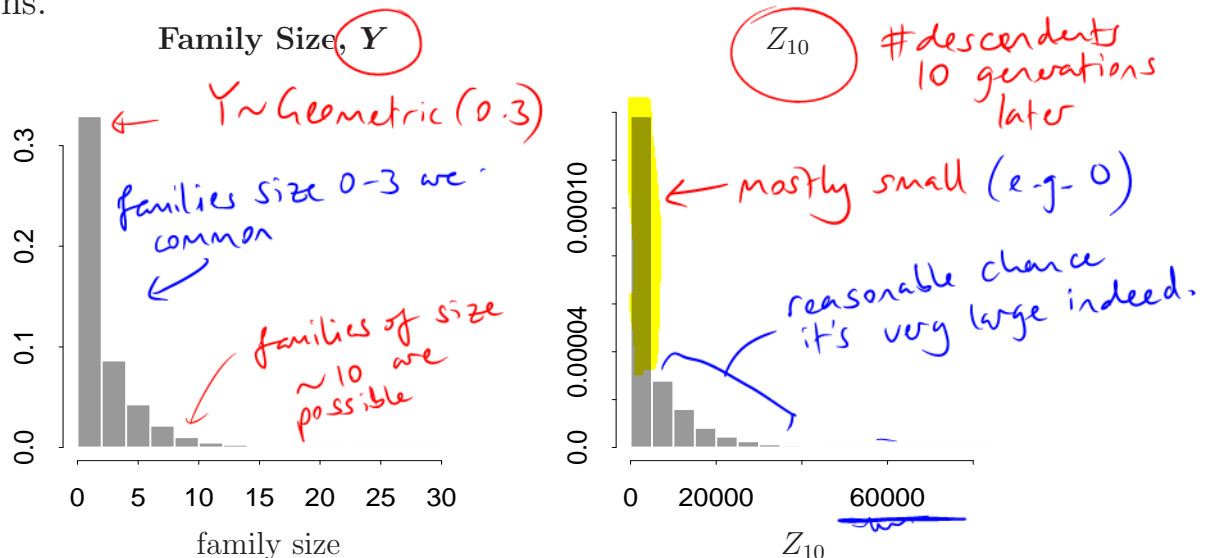
## Simulation 1: $Y \sim \text{Geometric}(p = 0.3)$

The following table shows the results from 10 simulations of a branching process, where the family size distribution is $Y \sim \text{Geometric}(p = 0.3)$. $\quad E(Y) = \frac{q}{p} = \frac{0.7}{0.3}$

$= 2 \cdot 3.$

| 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 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 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 |

*mean = 4617*

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

*= based on observation*

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.

**Family Size, $Y$**

$Z_{10}$ — #descendants 10 generations later



*$Y \sim \text{Geometric}(0.3)$*
*families size 0-3 are common*
*families of size ~10 are possible*
*mostly small (e.g. 0)*
*reasonable chance it's very large indeed.*

family size

$Z_{10}$

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

```
Mean of Zn:        4617.2
Variance of Zn:  53937785.7
```

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? No.

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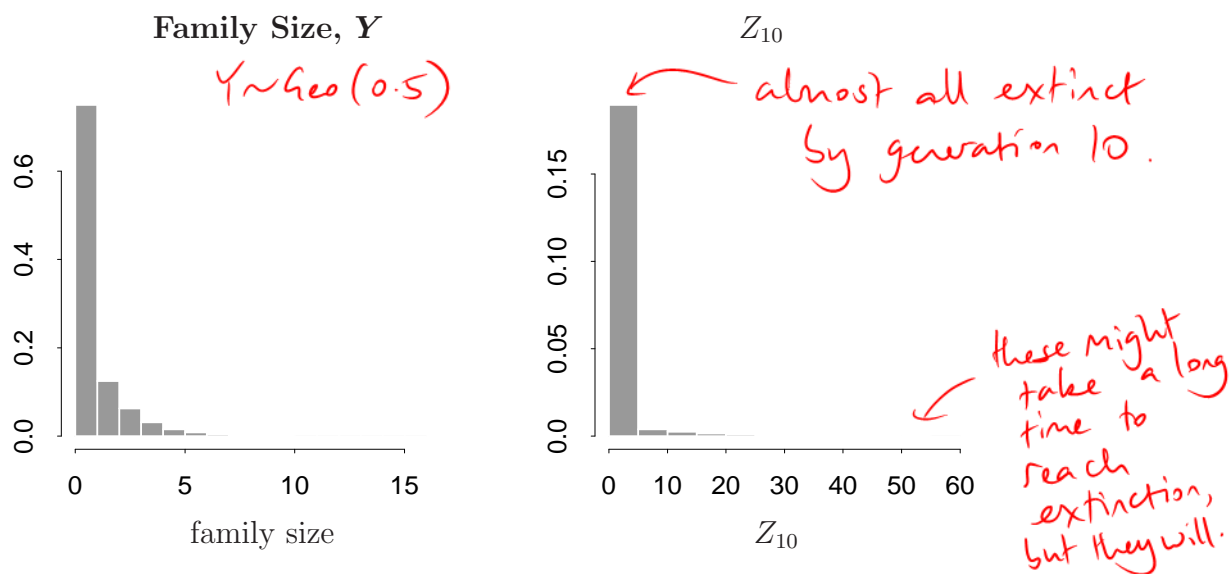
## Simulation 2: $Y \sim \text{Geometric}(p = 0.5)$

$\mathbb{E}(Y) = \dfrac{q}{p} = \dfrac{0.5}{0.5} = 1.$

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 $Y \sim \text{Geometric}(p = 0.5).$

| 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  |

This time, almost all the populations become extinct. We will see later that this value of $p$ (just) guarantees eventual extinction with probability 1. *because* $\mathbb{E}Y = 1$.

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.



**Family Size, $Y$** — *$Y \sim \text{Geo}(0.5)$*

$Z_{10}$ — *almost all extinct by generation 10.*

*these might take a long time to reach extinction, but they will.*

Proportion of samples extinct by generation 10:  0.9108

```
Summary of Zn:
     Min  1st Qu  Median    Mean  3rd Qu    Max
       0       0       0   0.965       0     56
```

*by simulation*

Mean of Zn:     0.965  ← $\mathbb{E}(Y) \sim 1$,    $\mathbb{E}(Z_{10}) \sim 1$.
Variance of Zn:  19.497  ← *variance is MUCH smaller but still greatly inflated (e.g. compare with Poisson, variance = mean).*

## What happens for larger values of $p$?

It was mentioned above that $Y \sim \text{Geometric}(p = 0.5)$ <u>just</u> 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.
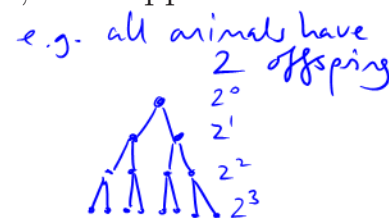
## 8.5 Mean and variance of $Z_n$

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 $\text{Var}(Y) = \sigma^2$. (Mean & variance of # offspring of a SINGLE individual.)

**Theorem 8.5:** Let $\{Z_0, Z_1, Z_2, \ldots\}$ 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

e.g. all animals have 2 offspring

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

**Proof:**

By p. 167, $Z_n = Y_1 + Y_2 + \cdots + Y_{Z_{n-1}}$ is a randomly stopped sum.

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

So by Section 3.4 (page 62),

$$\mathbb{E}(Z_n) = \mathbb{E}(Y)\,\mathbb{E}(Z_{n-1})$$
$$= \mu\,\mathbb{E}(Z_{n-1})$$
$$= \mu\{\mu\,\mathbb{E}(Z_{n-2})\}$$
$$= \mu^2\,\mathbb{E}(Z_{n-2})$$
$$\vdots$$
$$= \mu^{n-1}\,\underbrace{\mathbb{E}(Z_1)}_{=\mu}$$

$\left\{\right.$ prove formally by induction

$$= \mu^{n-1} * \mu$$
$$= \mu^n . \qquad \square$$

***Examples:*** Consider the simulations of Section 8.4.

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:

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

   The theoretical value, $4784$ compares well with the sample mean from 5000 simulations, $4617$ (p. 172).

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

   So $E(Z_{10}) = \mu^{10} = 1^{10} = 1.$

   Compares well with sample mean of $0.965$ (p. 173.)

---

## Variance of $Z_n$

**Theorem 8.5:** Let $\{Z_0, Z_1, Z_2, \ldots\}$ 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

Read if you like
(higher priorities in 2012)

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

**Proof:**

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

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

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

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

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

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

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

Also,

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

**Find $V_n$ by repeated substitution:**

$$V_1 \;=\; \sigma^2$$

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

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

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

$$\vdots \quad \text{etc.}$$

Completing the pattern,

$$V_n \;=\; \mu^{n-1} \sigma^2 \left(1 + \mu + \mu^2 + \ldots + \mu^{n-1}\right)$$

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

$$=\; \mu^{n-1} \sigma^2 \left(\frac{1 - \mu^n}{1 - \mu}\right). \qquad \text{Valid for } \mu \neq 1.$$

(sum of first $n$ terms of Geometric series)

When $\mu = 1$ :

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

Hence the result:

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

---

**Examples:** Again consider the simulations of Section 8.4.

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

$$\sigma^2 = \text{Var}(Y) = \frac{q}{p^2} = \frac{0.7}{(0.3)^2} = 7.78.$$

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

( compare with $5.39 \times 10^7$ from the simulation ).

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

$$\sigma^2 = \text{Var}(Y) = \frac{q}{p^2} = \frac{0.5}{(0.5)^2} = 2.$$

Using formula for $\text{Var}(Z_n)$ when $\mu = 1$, we get

$$\text{Var}(Z_{10}) = \sigma^2 n = 2 * 10 = 20.$$

( compare with $19.5$ by simulation.)