# **Cancer: A Disease of the DNA**

## Introduction

Cancer is a group of diseases in which cells grow and spread unrestrained throughout the body. Cancers can arise in nearly any type of cell that retains the ability to divide. Although there are more than 100 forms of cancer, the basic processes underlying all of them are very similar. The process by which normal cells become cancerous is called *carcinogenesis*.

Cancers stem from mistakes and misapplication of cellular mechanisms, the cell's inability to heed normal growth and division controls or to undergo self-destruction, called apoptosis, when it detects that it is damaged. Normal cells are part of a cellular community and coordinate their activities with those of their neighbors especially regarding growth and division. Cancerous cells ignore cellular controls and even produce false signals for coercing their neighbors to help them. This errant behavior comes about due to the accumulation of small mutations, changes to the cellular genome that are perpetuated in cell reproduction. Two gene classes play major roles in choreographing the cellular life cycle: proto-oncogenes initiate cell growth and division, and tumor suppressor genes inhibit cell growth and division. When proto-oncogenes go awry and become oncogenes, they maintain continuous growth signals, like a car with the accelerator pedal stuck on. By contrast, dysfunctional tumor suppressor genes are like a car with no brakes. In order for a tumor to develop, mutations usually must occur in several genes.

# **12.1** Cell Growth and Division Is an Involved and Complicated Process

In Chapter 13, we consider the cell's growth and division cycle from a functional point of view. Here we must take a close look at the biochemistry of the sequence of steps leading to the creation of a new cell in some detail. We need to know the proteins involved and to understand their roles in the cycle.

## A cell cycle involves the orchestrated activity of intracellular enzymes.

Each cell of the body lives as a member of a community of cells. The cell's community are the nearby cells of its particular tissue type. Together the community jointly controls the growth and division of its members. Cell division begins within a cell when it receives growth stimulatory signals transmitted by other cells of the community. These signals are in the form of protein *growth factors*, which move in the interstitial space between cells and bind to specific receptors on the surface of the cell. The receptors in turn signal proteins within the cell. The process is complicated in that several proteins of the cytoplasm are involved in a chain, or *stimulatory pathway*, of signaling, which ultimately ends in the nucleus of the cell.

In counterpoint to the stimulatory pathway, there is also an *inhibitory pathway*. Its signal transmitted to the cell nucleus is to hold off cell division.

The target of this activity in the nucleus is the *cell cycle clock*, also discussed in Section 13.1. The clock integrates all the stimulatory and inhibitory signals and decides whether the cell should undergo a division cycle, or *cell cycle*.

The cell cycle consists of four stages. To divide, the cell must double its genome; this stage is called the S, or *synthesis*, *phase*. Afterward, it must halve that genome in the *mitosis*, or M, phase. Between M and S is the *first gap phase*, or  $G_1$ , and between S and M is the *second gap phase*, or  $G_2$ . In  $G_1$ , the cell increases in size and at some point decides whether to divide. In  $G_2$ , the cell undergoes the necessary preparations for cell division, and in M the cell divides in half, sending a full complement of its chromosomes to each daughter cell. Each daughter cell immediately enters stage  $G_1$ but normally does not proceed to S right away. Instead, it then enters a resting, or  $G_0$ , stage. Later it will reenter  $G_1$  and decide to proceed on to S or return to  $G_0$ . The term *interphase* refers to any stage of the cell cycle other than M.

The cell cycle is quite complicated and is mediated at each step by a variety of molecules; principal among them are proteins called *cyclins* and *cyclin-dependent kinases* (CDKs). See Figure 12.1.1.

Cyclin D and CDK4 orchestrate  $G_1$ . Starting with  $G_1$ , a rising level of cyclin D binds to CDK4 and signals the cell to prepare the chromosomes for replication. It achieves this effect by nullifying the effect of the protein pRB, which exerts powerful growth-inhibitory control.

Cyclins A and E and CDK2 mediate phase S. A rising level of *S*-phase promoting factor (SPF), cyclin A bound to CDK2, enters the nucleus and prepares the cell to duplicate its DNA and its centrosomes. As DNA replication proceeds, cyclin E is destroyed, and the level of mitotic cyclins A and B begins to rise (in  $G_2$ ).

In stage M, a complex of cyclins A and B and CDK1, called *M-phase promoting factor*, initiates several things: the assembly of the mitotic spindle, the breakdown of the nuclear envelope, and the condensation of the chromosomes.

These events take the cell to the metaphase of mitosis. At this point, the M-phase promoting factor activates the *anaphase-promoting complex* (APC), which

• allows the sister chromatids at the metaphase plate to separate and move to the poles (this is *anaphase*), completing mitosis;



**Fig. 12.1.1.** Biochemistry of the cell cycle. An arrowhead path indicates a promoting signal; a blunthead path indicates an inhibiting signal.

- destroys cyclin B—it does this by attaching it to the protein *ubiquitin*, which targets it for destruction by proteasomes;
- turns on the synthesis of cyclin D for the next turn of the cycle;
- degrades *geminin*, a protein that has kept the freshly synthesized DNA in S phase from being rereplicated before mitosis.

In addition to these mechanisms for conducting the cell cycle, there are also checking mechanisms. There is a check of the DNA for errors in  $G_1$  before the cell enters S phase. There is also a check that all the DNA has been copied before the cell can leave S phase. The principal protein involved in this is p53. The gene for p53 is recessive, since both copies must be defective for p53 to fail. The protein p53 is also an important mediator of *apoptosis*, leading defective cells to undergo self-destruction.

After S phase, in G<sub>2</sub>, there is again a DNA check for damage. A molecule involved in this is *ataxia telangiectasia mutated*, or ATM.

There are also checkpoints in M phase. In one example, a check is made that the spindle fibers are properly attached. A central molecule in this check is *mitotic arrest deficient* protein, or MAD.

# 12.2 Types and Stages of Cancer

Cancers that arise from cells that cover external or internal body surfaces are *carcinomas*. The most common examples of these are lung, breast, and colon cancer. *Sarcomas* are cancers that arise from cells in supporting tissues such as bone, cartilage, fat, connective tissue, and muscle. *Lymphomas* arise in the lymph nodes and tissues of the immune system, and *leukemias* are cancers of the immature blood cells that develop in the bone marrow.

An increase in the number of dividing cells creates a growing mass of tissue called a *tumor* or *neoplasm*. There are fast- and slow-growing tumors. As in any tissue, in order to continue growing, tumors must have access to a blood supply. Blood vessels may be attracted to a growing tumor in a process called *angiogenesis*.

Tumor cells may become capable of invading nearby tissues—this is called *invasion*—or, in its extreme form, of invading distant tissues by a process called *metastasis*. In metastasis, tumor cells penetrate a lymphatic or blood vessel, into which they shed dividing cells. There the tumor cells face the body's normal immune response, which kills a large number of them. But any that escape can reach and invade some hitherto unaffected tissue.

A *benign* tumor is one that cannot spread by invasion or mestasts and remains local. By contrast, a *malignant* tumor is capable of both; generally the term cancer refers to a malignant tumor.

There are several stages to cancer development. *Hyperplasia* is tissue growth due to an excessive cell division rate in which cell structure and the arrangement of cells are normal. Hyperplasia can be due to an irritating stimulus. *Dysplasia* is excessive cell proliferation along with the loss of normal tissue arrangement and normal cell structure to some degree. Sever dysplasia is referred to as a *carcinoma in situ* and is uncontrolled cell growth. Nevertheless, at this stage, the tumor is not invasive or metastatic.

In addition, tumors are graded for degree of progression. If the cells are highly abnormal in appearance and there is a large number of dividing cells, the grade is III or IV. Cancers less progressed have a grade of I or II.

Cancer, being a disease of the DNA, is amenable to study from the standpoint of genetics. In Section 14.4, we delve into this topic briefly and note some of the methods being tried to bring the disease under control.

## 12.3 The Role of Proto-Oncogenes and Tumor Suppressor Genes

As mentioned above, proto-oncogenes are responsible for the stimulatory pathways leading to cell division. When a proto-oncogene mutates, leading to an overexcitation of its function in the pathway, the mutated gene is termed an *oncogene*. All

stages of the growth and division pathway can be subverted starting with the receptor proteins embedded in the cell surface. For example, the overactive receptor Erb-B2 is associated with breast cancer.

The next step in the pathway is the signal cascade operating in the cytoplasm. Here, too, oncogenes are responsible for overactive proteins. One example is the *ras* family. Proteins coded for by mutant *ras* genes are continuously stimulatory even in the absence of active growth factor receptors. Hyperactive *ras* proteins are found in one-fourth of all human tumors, including carcinomas of the colon, pancreas, and lung.

Oncogenes are implicated in growth and division pathologies in the nucleus as well. An example of this relates to the *myc* family of genes, coding for altered activity of various transcription factors. Such *myc* proteins are constantly at high levels despite the absence of their usual antecedents in the stimulatory pathway.

## Tumor suppressor genes are generally recessive.

Equally troublesome, aberrant *tumor suppressor genes* (TSG) can be at fault. Normal TSGs produce proteins that restrain cell growth and division, but when a TSG is defective, this control is absent. Defective tumor suppressor genes differ from oncogenes in the important respect that they are recessive, and hence both genes of the chromosomal pair must be defective for the encoded protein to be dysfunctional.

In some cases, tumor suppressor proteins can hold the stimulatory effect of an oncogene in check, and cancer will develop only when some component of the suppressor system fails. In one example, the substance *transforming growth factor beta* (TGF- $\beta$ ) is a regulatory component in the growth of normal cells. However, some colon cancers overcome TGF- $\beta$  by inactivating the gene that codes for its surface receptor. In another example, this time in the cytoplasm, some cancers evade the effect of the protein produced by the DPC4 gene by inactivating that gene. The same goes for the p15 gene, which codes for a protein that works in concert with TGF- $\beta$ . Normally, p15 works by blocking the binding of cyclin D with CDK4, thus preventing the cell going from G<sub>1</sub> to S. Should its gene become inactive, the cell loses this control.

The suppressor gene NF-1 is an example of a control that acts directly on a specific stimulatory protein. The suppressor produced by the NF-1 gene operates by blocking the growth signals sent by the *ras* protein. When NF-1's protein is dysfunctional, the stimulatory signal goes unchecked.

#### p53 plays a central role in the health of the cell.

One of the most prominent suppressor proteins is p53, which is multiroled. It checks the health of the entire cell generally: It monitors the integrity of the chromosomal DNA, it monitors the successful completion of the cell cycle, and it is the key protein involved in apoptosis. Half of all types of human cancers lack a functional p53 protein.

Regarding the effect of TSGs, it has been shown that the external introduction of a missing tumor suppressor protein can restore the control that that protein exhibits

and hence arrest an incipient cancer. Thus it may be possible to treat cancers that arise from dysfunctional TSGs pharmacologically.

## Cell death and cell division limits are also subverted.

If a cell's essential components become damaged, including the chromosomal DNA, then the cell will initiate events that lead to its own destruction; this is called *apoptosis*. As previously mentioned, the p53 protein generally monitors the health of the cell and triggers apoptosis if warranted. Among the DNA damage that can be detected is the conversion of a proto-oncogene to an oncogene or damage to a TSG.

Of course, if the gene coding for p53 itself becomes damaged, then apoptosis may not be possible for the cell. In addition, cancer cells devise ways of evading apoptosis. One of these is through mechanisms that inactivate the p53 encoding gene. Another is to make excessive amounts of the protein Bcl-2, which counters p53.

It was once thought that radiation therapy disrupted a cell so badly that it died. But now it is known that, in fact, only minor damage is inflicted on the cell, and it is apoptosis that kills the cell. The implication is that if p53 is absent or dysfunctional, then radiation therapy is far less effective.

A second mechanism that normally operates to kill cancer cells relates to the fact that cells may divide only a limited number of times, up to 50 to 60 depending on the type of cell. At that point the cell is said to be *senescent*; see also Section 5.1.

The mechanism for this has to do with segments on the ends of the chromosomal DNA strands known as *telomeres*. Each time a cell's chromosomes are replicated, the telomeres shorten a bit. Once the telomere length falls below a threshold, a signal goes out instructing the cell to enter senescence. If chromosomal duplication should proceed despite this signal, chromosomes fuse with one another, causing the cell to die.

Replication of a chromosome is accomplished by a protein that moves over the length of its surface and reads the base pairs. By way of analogy, think of a railroad car that is pulled over the rails to clean the track. When the cleaning train reaches the end of the line, the track is not cleaned to the very end because the locomotive is in the way. In the same way, the tip of the chromosomes cannot be read and duplicated, hence the function of the telomeres. The unread tip is the length by which the telomere shortens on each replication.

The proteins pRB and p53 are involved in detecting telomere length and signaling senescence. But cells that have mutations to the genes for either of these proteins will continue to divide anyway. Thus as explained above, eventually the telomeres disappear altogether and the cell dies. However, there is a way that a cell can circumvent this fate and become immortal.

The enzyme telomerase systematically replaces lost telomere segments, thereby enabling the chromosomes to be replicated endlessly. This enzyme is absent in most normal cells but present in almost all tumor cells. Thus not only does the tumor cell not die, but there is a second pathological effect. Immortality gives the cell time to accumulate other mutations.

Yet another way the necessary mutations can come about is through nonfunctional DNA repair proteins. The cell possesses proteins that move over the chromosomal

DNA and detect and repair copying mistakes and other abnormalities. As mentioned above, one of these is ATM. But these proteins may become dysfunctional, for example, by means of a mutation to their own genes. In this case, unrepaired mutations within the cell can accumulate very quickly. ATM is implicated in 10% of inherited breast cancers.

## It usually takes several mutations or "hits" for malignancy.

Most malignancies are the product of tumor cells that have acquired several mutations. These may include several or all of the following:

- an overstimulated cell cycle pathway,
- a defective tumor suppressor pathway,
- the ability to circumvent apoptosis,
- the ability to produce telomerase,
- a mechanism for attracting blood vessels (angiogenesis),
- the ability to invade nearby tissues (invasion),
- the ability to metastasize.

Normally, it takes decades for an incipient tumor to collect enough mutations for malignant growth, but there is a mechanism that greatly accelerates the rate—an inherited mutation. If a germ cell harbors some mutation, then that mutation is part of every cell of the body. For example, suppose a germ cell has a defective copy of the TSG gene. Since these genes are recessive, the defect will not manifest itself so long as its twin on the homologous chromosome remains intact. But now every cell of the body is a target. Should a mutation occur to the remaining working gene in any one of them, the tumor suppression mechanism will now fail for that cell. Such an event is vastly more probable than accumulating two defects by chance to the same cell. We investigate this reality mathematically through a retinoblastoma study in Section 12.4.

Although inherited genetic alterations are a factor in some cases, nevertheless 80–90% of cancers occur in people with no family history of the disease.

The progression of cells to cancer is reminiscent of Darwinain evolution. Cells that acquire any of these mutations have the capacity for reproducing themselves to a greater extent than normal cells. This has the consequence that their kind outgrow neighboring cells with several detrimental consequences. The mildest of these is that these aberrant cells no longer participate fully in the intended function of the tissue while at the same time garnering more and more of the tissue's resources such as space and nutrients. Thus the normal cells of the tissue are at a competitive disadvantage. More seriously, they beget a colony of cells capable of further mutation.

## 12.4 A Model for Retinoblastoma

Retinoblastoma is a cancer of the retina and is one of the simplest cancers in terms of mechanism. It results when a dividing retinal cell lacks a tumor suppressor protein

as a result of defects in the corresponding gene to both homologous chromosomes. Retinal cells grow and divide from early in fetal development and continue growing after birth until about age 2 or 3 years. Thus retinoblastoma is a childhood disease.

One way a retinal cell can suffer the two defects, or *hits*, to its tumor suppressor gene is by chance alone. Such a case is called *somatic* and is thought to account for 60% of all cases of the disease. Actually, this is a very rare occurrence with an incidence rate in the general population of about 30 per million.

The other possibility is for an individual to inherit a defective gene for the disease; such a person is called a *carrier*. In this event, every cell of the body, including all retinal cells, has one hit already. Consequently, if a mutation now occurs to any retinal cell at all during the growth phase, that cell will have two hits, and tumor growth ensues. These are called *germinal* cases and account for the balance of the cases, about 40%.

#### Patient data.

In his study of the disease, Knudson surveys cases of retinoblastoma and presents the data shown in Table 12.4.1. In fact, this purely mathematical study suggested the existence of tumor suppressor proteins, which led to their discovery. As reported in the table, retinoblastoma cases are observed as to sex, onset age, laterality (that is, which eye holds tumors), number of tumors, and family history of the disease. In addition, Knudson further presents the summary data of Table 12.4.2. Among this data is an estimate of the fraction of the cases that are germinal and, by complement, somatic. This estimate is made using the assumption that somatic cases have only one tumor, due to the unlikelihood of two hits; cf. statistic 13. A key item researched by Knudson is "carriers never affected," statistic 1, which can be used to determine the gene mutation rate p, as we will see below. Note that Knudson did not assess statistic 6, "carriers among the population," but we will be able to do so from our analysis.

We formulate a simple model to account for the data, namely that a retinal cell grows and matures for a period of time, after which it divides. In the process of division, a mutation randomly occurs to the tumor suppressor gene with probability p. We assume that there is no cell death. The two daughter cells are genetically the same as their shared parent, except that with probability p, one of them has mutated. (Probability p is so small that we ignore the possibility that both daughter cells mutate.) A mutant cell then begets a clone of like cells.

To further simplify the model, we assume that all cells wait the same period of time, one unit, and hence we have a synchronous population. A more advanced model allows for random growing periods and therefore overlapping generations, but the results are essentially the same.

Our first task is to determine how many such cell divisions, T, retinal cells undergo. It is known that there are about 2 million retinal cells per eye, or  $4 \times 10^6$  counting both eyes, all starting from a single cell. Therefore, there must be T = 22 division cycles, since  $2^{22} = 4,194,304$ . As mentioned above, this takes about two to three years on average. This result can also be derived by solving a simple recurrence

(a) Bilateral cases						(b) Unilateral cases				
		Age at	Nur	nber of				Age at		
		diagnosis	tu	mors	Family			diagnosis	Number of	Family
Case	Sex	(months)	left	right	history	Case	Sex	(months)	tumors	history
1	F	8	*	*	no	24	F	48	*	no
2	Μ	3	*	*	no	25	Μ	22	*	no
3	F	11	*	*	no	26	Μ	33	*	no
4	Μ	2	*	1	no	27	Μ	38	*	no
5	Μ	60	1	*	affected sib.	28	F	47	*	no
6	Μ	22	*	*	no	29	Μ	50	*	no
7	F	4	3	*	no	30	Μ	32	*	no
8	F	18	2	*	no	31	Μ	28	*	no
9	F	30	*	*	no	32	F	31	*	no
10	Μ	3	2	*	no	33	F	29	*	no
11	F	6	*	1	no	34	F	21	*	no
12	Μ	7	*	2	affected sib.	35	Μ	46	*	no
13	Μ	9	3	*	no	36	F	36	*	no
14	F	4	5	*	no	37	F	73	*	no
15	F	13	*	*	no	38	Μ	29	*	no
16	Μ	18	*	4	no	39	F	15	*	no
17	F	24	*	*	no	40	Μ	52	*	no
18	F	44	1	*	no	41	Μ	24	*	no
19	F	5	*	*	no	42	Μ	8	*	no
20	Μ	12	*	1	no	43	F	19	*	no
21	Μ	3	*	*	no	44	Μ	36	*	no
22	Μ	12	*	1	no	45	F	34	*	no
23	Μ	15	1	*	father	46	F	27	*	no
						47	Μ	10	*	no
						48	F	8	*	no

Table 12.4.1. Cases of retinoblastoma (\* indicates no information).

relation (cf. Section 2.5). Let N(t) be the number of retinal cells after t cell divisions. Then N(0) = 1, and for any t = 1, 2, ...,

$$N(t) = 2N(t-1).$$
(12.4.1)

The solution is

 $N(t) = 2^t, \quad t = 0, 1, 2, \ldots;$ 

try it!

## Germinal cases determine the mutation probability p.

Consider a germinal case. Should a mutation occur to some retinal cell of a carrier, the result is a doubly mutated cell and hence a tumor. We are not interested in the number of such doubly mutated cells per se but rather in the number of tumors, which here is the same as the number of mutation events occurring to the growing tissue.

	Empirical	Calculated using
Statistic	value	$p = 7.14 \times 10^{-7}$
1. carriers never affected, used for $p$	1-10%	0.05
2. average number of tumors among germinal cases	3	2.99
3. unilateral cases among germinal cases, $g_1$	25-40%	0.337
4. bilateral cases among germinal cases	60–75%	0.663
5. bilateral cases among all cases, <i>B</i>	25-30%	0.27
6. carriers among general population, $f$		27 per million
7. incident rate of somatic cases, <i>u</i>	30 per million	40 per million
8. germinal cases among all cases	35-45%	0.387
9. somatic cases among all cases	55-65%	0.613
10. unilateral cases among all cases	70–75%	0.750
11. unilateral and hereditary cases among all cases	10-15%	0.137
12. germinal cases among unilateral cases	15-20%	0.183
13. unilateral cases among nonhereditary cases	100%	assumption

Table 12.4.2. Summary data.

Let X(t) be the random variable denoting the number of tumors at time t, and let  $x_k(t)$  be the probability that there are k tumors at time t, that is,  $x_k(t) = \Pr(X(t) = k)$ , k = 0, 1, 2, ...

A powerful tool in the derivation of equations for the  $x_k$  is the *probability-generating function* G(s, t), a version of which will also be prominent in Section 13.7. In this case, G(s, t) is defined as

$$G(s,t) = \sum_{0}^{\infty} x_k(t) s^k.$$
 (12.4.2)

Evidently, *G* is a function of two variables, *s* and *t*, a polynomial in *s* whose coefficients are the  $x_k(t)$ . The variable *s* merely acts as a bookkeeping device here in that it manages the  $x_k$ , but it is a most useful bookkeeping device indeed. Even so, treating *s* as a variable and setting s = 1, we get

$$G(1,t) = x_0(t) + x_1(t) + x_2(t) + \dots = 1, \qquad (12.4.3)$$

which sums to 1 because the xs account for all possibilities at time t: Either there is none, or 1, or 2, and so on. Also, we can recover all the  $x_k$  from G by differentiating with respect to s and substituting s = 0. Thus

$$G(0,t) = x_0(t)$$

gives  $x_0(t)$ , and the partial derivative

$$\frac{\partial G}{\partial s} = x_1(t) + 2x_2(t)s + 3x_3(t)s^2 + \cdots,$$

along with the substitution s = 0, gives  $x_1(t)$ ,

$$\left. \frac{\partial G}{\partial s} \right|_{s=0} = x_1(t).$$

Continuing to differentiate with respect to s and substituting s = 0, one sees that

$$x_k(t) = \left. \frac{1}{k!} \frac{\partial^k G}{\partial s^k} \right|_{s=0}.$$
 (12.4.4)

In addition to the  $x_k$ , we can also calculate expectations from G(s, t). The expectation E(X(t)), or expected value, of X(t) is defined as the sum over all its possible values weighted by the probability of each, so

$$E(X(t)) = 1 \cdot x_1(t) + 2 \cdot x_2(t) + 3 \cdot x_3(t) + \cdots$$

But this is exactly the partial derivative of G with respect to s evaluated at s = 1,

$$\left. \frac{\partial G}{\partial s} \right|_{s=1} = x_1(t) + 2sx_2(t) + 3s^2x_3(t) + \dots + |_{s=1} = \sum_{k=1}^{\infty} kx_k(t).$$
(12.4.5)

One more observation on G. Its square is also an infinite series in s, whose coefficients take a special form,

$$G^{2} = (x_{0} + x_{1}s + x_{2}s^{2} + \dots)(x_{0} + x_{1}s + x_{2}s^{2} + \dots)$$
  
=  $x_{0}^{2} + (x_{0}x_{1} + x_{0}x_{1})s + (x_{0}x_{2} + x_{1}^{2} + x_{2}x_{0})s^{2} + \dots$   
=  $\sum_{k=0}^{\infty} \left(\sum_{k_{1}+k_{2}=k} x_{k_{1}}x_{k_{2}}\right)s^{k}.$  (12.4.6)

With these preliminaries behind us, we set our sights on finding  $x_0(t)$ . Remarkably, it is easier to calculate all the  $x_k$  than just  $x_0$ ! Start by decomposing on the possibilities at t = 1; either there is a mutation on this first cell division or not. "Given"—symbolized by |—there is a mutation, the probability that there will be k mutations by time t is exactly the probability that there are k - 1 in the remaining time, t - 1, so the contribution here is

$$Pr(k \text{ tumors at } t \mid 1 \text{ at } t = 1) Pr(a \text{ mutation at } t = 1) = x_{k-1}(t-1)p.$$

On the other hand, if there is no mutation at time t = 1, then there must be k mutations in the remaining time t - 1 stemming from one or both of the two daughter cells. If one of the daughter cells leads to  $k_1$  tumors in the remaining time, where  $k_1$  could be any of 0, 1, 2, ..., the other must lead to  $k_2 = k - k_1$ . Thus the contribution here is

$$Pr(k \text{ tumors at } t \mid \text{none at } 1) Pr(\text{no mutation at } t = 1)$$

$$= \left(\sum_{k_1+k_2=k} x_{k_1}(t-1)x_{k_2}(t-1)\right)(1-p).$$

Combining these two contributions, we have the equation

$$x_{k}(t) = \Pr(k \text{ tumors at } t \mid \text{none at } 1) \Pr(\text{no mutation at } t = 1) + \Pr(k \text{ tumors at } t \mid 1 \text{ at } t = 1) \Pr(\text{mutation at } t = 1) = (1 - p) \sum_{k_{1}+k_{2}=k} x_{k_{1}}(t - 1)x_{k_{2}}(t - 1) + px_{k-1}(t - 1).$$
(12.4.7)

Multiply both sides by  $s^k$  and sum over k to get

$$\sum_{k=0}^{\infty} x_k(t) s^k = (1-p) \sum_{k=0}^{\infty} \left( \sum_{k_1+k_2=k} x_{k_1}(t-1) x_{k_2}(t-1) \right) s^k + ps \sum_{k=1}^{\infty} x_{k-1}(t-1) s^{k-1}.$$

It remains only to recall that the double sum on the right is exactly  $G^2(s, t-1)$ . Thus we arrive at the fundamental equation

$$G(s,t) = (1-p)G^{2}(s,t-1) + spG(s,t-1).$$
(12.4.8)

As noted above, from the generating function we can feasibly calculate many of the properties of interest. One of these is incidence expectation. The expected number of tumors, E(X(t)), is given by

$$E(X(t)) = \left. \frac{\partial G}{\partial s} \right|_{s=1}$$
  
= 2(1-p)G(1,t-1)  $\frac{\partial G(1,t-1)}{\partial s} + pG(1,t-1) + ps \frac{\partial G(1,t-1)}{\partial s}.$ 

Remembering that G(1, t) = 1 for all t gives the following recurrence relation and initial value (the first and last terms combine):

$$E(X(t)) = (2 - p)E(X(t - 1)) + p, \quad E(X(1)) = p.$$

This is easily solved to give

$$E(X(t)) = \frac{p}{1-p} \left( (2-p)^t - 1 \right).$$
(12.4.9)

We may also calculate some individual probabilities.

Since  $x_0(t) = G(0, t)$ , substituting s = 0 in (12.4.8) gives the following recurrence relation and initial value:

$$x_0(t) = (1-p)x_0^2(t-1), \quad x_0(1) = 1-p.$$

\_

This, too, is easily solved, and yields the probability that no tumor will occur by time t,

$$x_0(t) = (1-p)^{2^t-1}.$$
 (12.4.10)

In a similar fashion, one obtains the probability that there will be one tumor,  $x_1(t)$ . From the definition of the generating function,  $x_1(t) = \frac{\partial G}{\partial s}|_{s=0}$ . Differentiating (12.4.8), setting s = 0, and using (12.4.10) gives

$$x_1(t) = 2(1-p)x_0(t-1)x_1(t-1) + px_0(t-1)$$
  
=  $(1-p)^{2^{t-1}-1} [2(1-p)x_1(t-1) + p], \quad x_1(1) = p.$  (12.4.11)

While this equation and, to a greater extent, those for the larger values of k are hard to solve in closed form (and we won't do that), they present no difficulty numerically. The reader can check that the next two equations for the  $x_k$  are given by

$$2x_2(t) = 2(1-p)[2x_0(t-1)x_2(t-1) + x_1^2(t-1)] + 2px_1(t-1), \quad x_2(1) = 0,$$
(12.4.12)

and

$$3!x_3(t) = 2(1-p)[3!x_1(t-1)x_2(t-1) + 3!x_0(t-1)x_3(t-1)] + 3!px_2(t-1), \quad x_3(1) = 0.$$
(12.4.13)

These may be calculated in MAPLE and MATLAB with the following programs:

```
MAPLE
> Digits:=60:
> x0:=array(1..22): x1:=array(1..22):
> x2:=array(1..22): x3:=array(1..22):
> expTumors:=(p,t)->(p/(1-p))*((2-p)^t-1):
> x0calc:=proc(p)
  global x0; local i;
    for i from 1 to 22 do
     x0[i]:= (1-p)^(2^i-1);
   od:
 end:
> x1calc:=proc(p)
  global x0, x1; local i;
   x1[1]:=p;
    for i from 2 to 22 do
     x1[i]:=x0[i-1]*(2*(1-p)*x1[i-1]+p);
   od:
  end:
> x2calc:=proc(p)
  global x0, x1, x2; local i;
   x2[1]:=0;
    for i from 2 to 22 do
     x2[i]:=(1-p)*(2*x0[i-1]*x2[i-1]+x1[i-1]^2)+p*x1[i-1];
   od.
 end:
> x3calc:=proc(p)
  global x0, x1, x2, x3; local i;
    x3[1]:=0:
    for i from 2 to 22 do
     x3[i]:=2*(1-p)*(x1[i-1]*x2[i-1]+x0[i-1]*x3[i-1])+p*x2[i-1];
    od;
 end:
> p:=0.0000007;
> x0calc(p); x1calc(p); x2calc(p); x3calc(p);
```

```
MATLAB
  % make an m-file, retinox0.m, with
  % function x0=retinox0(p)
  % for i=1:22
  % x0(i)=(1-p)^(2^i-1);
  % end
  % make an m-file, retinox1.m, with
 % function x1=retinox1(p,x0)
  % x1(1)=p;
  % for i=2:22
  %
     x1(i)=x0(i-1)*(2*(1-p)*x1(i-1)+p);
 % end
  % make an m-file, retinox2.m, with
  % function x2=retinox2(p,x0,x1)
  % x2(1)=0:
  % for i=2:22
 %
      x2(i)=(1-p)^{*}(2^{*}x0(i-1)^{*}x2(i-1)+x1(i-1)^{2})+p^{*}x1(i-1);
 % end
  % make an m-file, retinox3.m, with
  % function x3=retinox3(p,x0,x1,x2)
  % x3(1)=0:
 % for i=2:22
 %
      x_{3(i)=2^{(1-p)^{(1-1)^{x}2(i-1)+x_{0(i-1)^{x}3(i-1))+p^{x}2(i-1)}}
 % end
> p=.0000007;
> x0 = retinox0(p)
> x1 = retinox1(p,x0)
> x2 = retinox2(p,x0,x1)
> x3=retinox3(p,x0,x1,x2)
```

## 12.5 Application to the Retinoblastoma Data

In this section, we apply the model to the summary statistics. The predicted results are reported on the right side of Table 12.4.2. It can be seen that the model is consistent with the data. First, we use statistic 1 to calculate p.

From (12.4.10) in conjunction with statistic 1, carriers never affected, which should be about 5%, we get a value for p,

$$0.05 = x_0(22) = (1-p)^{2^{22}-1}, \qquad \log(1-p) = \frac{\log(0.05)}{2^{22}-1} = -0.0000007142.$$

So (1-p) = 0.9999992858 and  $p = 7.14 \times 10^{-7}$ . Thus the probability of a mutation occurring upon cell division that renders a tumor suppressor gene dysfunctional is about  $7 \times 10^{-7}$ . A very small value, but in the production of a full complement of retinal cells, there are about  $4 \times 10^6$  cell divisions. (The last of the synchronous cell divisions involves  $2 \times 10^6$  cells by itself.)

With this value of p, (12.4.9) gives the expected number of tumors in the germinal case, statistic 2, to be 2.99; compare with 3.

Denote by  $g_1$  the unilateral cases among germinal cases, statistic 3. It is calculated by the infinite series

$$g_1 = x_1(22) + \frac{1}{2}x_2(22) + \frac{1}{2^2}x_3(22) + \dots$$
 (12.5.1)

for the following reason: There could be one tumor after all cell divisions; this is  $x_1(22)$ . If there are two tumors, with probability  $\frac{1}{2}$  both are in the same eye. Similarly, if there are three, with probability  $\frac{1}{4}$  all three are in the same eye because with probability  $\frac{1}{2}$  the second is in the same eye as the first and with another probability  $\frac{1}{2}$  the third is in that eye, too. And so on, but the remaining terms are small, so we stop with these three terms. Now step through the recurrence relations, (12.4.11), (12.4.12), (12.4.13), numerically with the value of *p* found above to get  $g_1 = 0.337$ ,

$$x_1(22) + \frac{1}{2}x_2(22) + \frac{1}{2^2}x_3(22) + \dots \approx 0.337.$$

Of course, statistic 4, bilateral cases among germinal cases, is the complement of this at 0.663.

#### The model also predicts carrier frequency.

Let f denote the fraction of carriers in the general population. Knowing f would allow us to estimate the other statistics of Table 12.4.2. Or we can use one of them to find f and then estimate the rest. The most reliable statistic is 5, the fraction of cases that are bilateral, about 27%; denote it by B:

$$B = \frac{\text{bilateral cases}}{\text{all cases}}.$$

Now by the assumption that all bilateral cases are germinal, statistic 13, the numerator will be

$$f(1-g_1).$$

For the denominator, the overall incidence rate of retinoblastoma, counting both germinal and somatic cases, is given by  $f(1 - x_0(22)) + (1 - f)u$ , where, as above, u is the probability of a somatic case. Hence we have for B,

$$B = \frac{f(1-g_1)}{f(1-x_0(22) + (1-f)u}.$$
(12.5.2)

Solve this for f to get

$$f = \frac{Bu}{1 - g_1 + Bu - B(1 - x_0(22))}.$$
 (12.5.3)

Substituting the numerical values, including  $u = 30 \times 10^{-6}$  from statistic 7, we get

$$f = 0.0000274,$$

or about 27 carriers per million in the general population.

Now the remaining statistics of Table 12.4.2 can be calculated. The fraction of germinal cases among all cases is given by

$$\frac{f(1-x_0(22))}{f(1-x_0(22)) + (1-f)u} = 0.387.$$
 (12.5.4)

Of course, the somatic cases are the complementary fraction.

The fraction of all cases that are unilateral is given by (explain)

$$\frac{f(x_1(22) + \frac{1}{2}x_2(22) + \frac{1}{2^2}x_3(22) + \dots) + (1 - f)u}{f(1 - x_0(22)) + (1 - f)u} = 0.750.$$
 (12.5.5)

The fraction of unilateral and heredity cases among all cases is given by

$$f(x_1(22) + \frac{1}{2}x_2(22) + \frac{1}{2^2}x_3(22) + \dots)f(1 - x_0(22)) + (1 - f)u = 0.137.$$
(12.5.6)

And finally, the fraction of unilateral cases that are germinal is

$$\frac{f(x_1(22) + \frac{1}{2}x_2(22) + \frac{1}{2^2}x_3(22) + \cdots)}{f(x_1(22) + \frac{1}{2}x_2(22) + \frac{1}{2^2}x_3(22) + \cdots) + (1 - f)u} = 0.183.$$
(12.5.7)

## 12.6 Persistence of Germinal Cases

It is natural to wonder how germinal cases arise: familial or new. In this section, we show that carriers are not persistent in society in that the genetic defect lasts at most two generations normally. Therefore, the condition is also the result of a chance mutation.

## Carriers can persist in the population.

Let  $q_k$  be the probability that a carrier zygote will survive and beget k offspring who are also carriers, and let  $F(s) = q_0 + q_1 s + q_2 s^2 + \cdots$  be the probability-generating function for the  $q_k$ . To compute the  $q_k$ , we also need the probabilities  $c_i$  that a surviving carrier will beget i offspring. Finally, let  $p_0$  be the probability that a carrier will survive to adulthood, taken as 0.05 from Table 12.4.2. Assuming that a carrier mates with a noncarrier, we get, for k > 0,

$$q_{k} = p_{0} \left( c_{k} \frac{1}{2^{k}} + c_{k+1} \binom{k+1}{k} \frac{1}{2^{k+1}} + c_{k+2} \binom{k+2}{k} \frac{1}{2^{k+2}} + \cdots \right),$$

since with probability  $\frac{1}{2}$  an offspring is a carrier or a noncarrier. And

$$q_0 = 1 - p_0 + p_0 \left( c_0 + c_1 \frac{1}{2} + c_2 \frac{1}{2^2} + \cdots \right).$$
 (12.6.1)

Standard theory provides that a trait will persist with probability 1 - V and die out with probability V, where V is the smallest fixed point of F, i.e., the smallest solution of F(L) = L; see also Section 13.7. Since already  $F(0) = q_0 > 0.95$ , L must be very close to 1.

Actually, L < 1 if and only if  $F'(1) \ge 1$ . But

$$F'(1) = q_1 + 2q_2 + 3q_3 + \cdots, \qquad (12.6.2)$$

and since each term is multiplied by  $p_0 = 0.05$ , even with extraordinary high values of the  $c_i$  for large i, F'(1) will be less than 1. For example, suppose  $c_5 = 1$ , i.e., an average carrier has five children. Then  $q_k = 0.05 \frac{\binom{k}{26}}{26}$ ,  $k = 1, \ldots, 6$  and  $F'(1) = \frac{192}{1280}$ .

## **Exercises/Experiments**

- 1. The first three statistics in Table 12.4.2 depend only on p. Plot these statistics as a function of p varying, say, from  $5 \times 10^{-7}$  to  $8 \times 10^{-7}$ . Recall that the code is given on p. 411 and  $g_1$  is given by (12.5.1). Explain your findings.
- 2. It can be shown that u, the probability of a somatic case, can also be calculated from p; in fact, u = 1 h(1, 22), where h(1, t) is given by the recurrence equation<sup>1</sup>

$$h(1,t) = 2px_0(t-1)h(1,t-1) + (1-2p)h^2(1,t-1).$$

As in Exercise 1 above, experiment with the behavior of u as a function of p. What value of p makes  $u = 30 \times 10^{-6}$ , 30 per million? What are the values of the first three statistics of Table 12.4.2 for this value of p?

- **3.** As in Exercise 2 above, *u* can be calculated knowing *p*. If also *B* is given a value, statistic 5 of Table 12.4.2, then all the other statistics can be calculated from these two; see (12.5.3), (12.5.4), (12.5.5), (12.5.6), and (12.5.7). For each of the four extreme values of statistics 1 and 5 of Table 12.4.2, namely,
  - (a) carriers never affected = 1%, B = 25%,
  - (b) carriers never affected = 1%, B = 30%,
  - (c) carriers never affected = 10%, B = 25%,
  - (d) carriers never affected = 10%, B = 30%,

tabulate the other statistics of the table.

- 4. It is not necessarily the case that the retinal cell divisions proceed uniformly in time. Suppose the time t in months from conception for the *i*th cell division,  $0 \le i \le 22$ , is given by  $i = 22(1 e^{-t/6})$ .
  - (a) How many cell divisions have occurred by birth, t = 9?

 $r_{n,0}(t) = \Pr(n \text{ with single mutation, no tumor at } t \mid a \text{ mutation at } 1)2p$ 

+  $Pr(n \text{ with single mutation, no tumor at } t \mid no mutation at 1)(1 - 2p)$ 

$$=2px_0(t-1)r_{n-2^{t-1},0}(t-1)+(1-2p)\sum_{n_1+n_2=n}r_{n_1,0}(t-1)r_{n_2,0}(t-1).$$

Define the generating function  $h(s, t) = \sum_{n} r_{n,0}(t)s^{n}$  and proceed as in the other derivations with generating functions.

<sup>&</sup>lt;sup>1</sup> Let  $r_{n,0}(t) = \Pr(n \text{ cells with a single mutation but no tumors at time } t)$ . Decompose on the possibilities at t = 1,

- 416 12 Cancer: A Disease of the DNA
  - (b) What is the chance the infant is born with a tumor if  $p = 7.14 \times 10^{-7}$  and what is the expected number of tumors at that time? (*Hint*: In the computer codes, stop the for loops at the value of *i* of part (a).)
  - (c) When does the 22nd cell division occur?
  - 5. Referring to Section 12.6, assume that each carrier mates with a noncarrier and has exactly five children. What is the minimum value of  $p_0$  that a carrier will survive to adulthood, in order that  $F'(1) \ge 1$  (see (12.6.2)), i.e., that the single tumor suppressor gene mutation persists?

## **Questions for Thought and Discussion**

- 1. In what ways will knowing the human genome be a help in fighting cancer?
- **2.** The retinoblastoma model measures time in terms of cell division generation. What observations are necessary to map this sense of time to real time?
- **3.** The mechanism in the text offers an explanation as to why the telomeres are shorter on the duplicate chromosomes. What about the original chromosomes used as the templates?

# **References and Suggested Further Reading**

- [1] B. Vogelstein and K. W. Kinzler, The multistep nature of cancer, *Trends Genet.*, **9**-4 (1993), 138–141.
- [2] A. G. Knudson, Jr., Mutation and cancer: Statistical study of retinoblastoma, Proc. Nat. Acad. Sci. USA, 68-4 (1971), 820–823.
- [3] A. G. Knudson, Jr., H. W. Hethcote, and B. W. Brown, Mutation and childhood cancer: A probabilistic model for the incidence of retinoblastoma, *Proc. Nat. Acad. Sci. USA*, 72-12 (1975), 5116–5120.
- [4] H. W. Hethcote and A. G. Knudson, Jr., Model for the incidence of embryonal cancers: Application to retinoblastoma, *Proc. Nat. Acad. Sci. USA*, **75**-5 (1978), 2453–2457.
- [5] G. M. Cooper, Oncogenes, 2nd ed., Jones and Bartlett, Boston, 1995.