

Chapter Title: Scaling Laws of Life, the Internet, and Social Networks

Book Title: Topics in Mathematical Modeling

Book Author(s): TUNG K. K.

Published by: Princeton University Press. (2007)

Stable URL: <https://www.jstor.org/stable/j.ctt1bw1hh8.5>

---

JSTOR is a not-for-profit service that helps scholars, researchers, and students discover, use, and build upon a wide range of content in a trusted digital archive. We use information technology and tools to increase productivity and facilitate new forms of scholarship. For more information about JSTOR, please contact support@jstor.org.

Your use of the JSTOR archive indicates your acceptance of the Terms & Conditions of Use, available at

<https://about.jstor.org/terms>



JSTOR

Princeton University Press is collaborating with JSTOR to digitize, preserve and extend access to *Topics in Mathematical Modeling*

# 2

## Scaling Laws of Life, the Internet, and Social Networks

### Mathematics required:

logarithms, log-log plots, geometric series, high school probability, simple binomial expansion

### Mathematics introduced:

self-similarity; first-order difference equation with variable coefficient and its solution by iteration

### 2.1 Introduction

Biology, the study of life forms, has historically been an empirical field, relying mainly on observations and classifications, but progress in the last few decades is beginning to transform it rapidly from a qualitative to a quantitative science. Mathematical modeling is playing an important role in this process. In this chapter we discuss the mathematics of networks, examples of which are the vascular blood network, which brings oxygen to the tissues served by the capillaries, and the plant vascular network, which transports nutrients from roots to leaves.

The World Wide Web has in recent decades been developed exponentially in size and complexity by random individuals attaching their web pages to sites of their choice, apparently without an overall design. Yet, when examined, the huge creature created in this way appears to contain structures in common with biological organisms. Advances in network theory have now allowed a study of these and other very complex networks, including social networks of people (such as actors) and the network formed by author citations. Hopefully these results will one day help us understand the even more complex neural networks in our brain.

### 2.2 Law of Quarter Powers

An important empirical law in biology is the allometric scaling law, which tells how an animal's property scales with its size or mass. (*Allo* comes from the Greek word *allos*, meaning "other.") So allometric

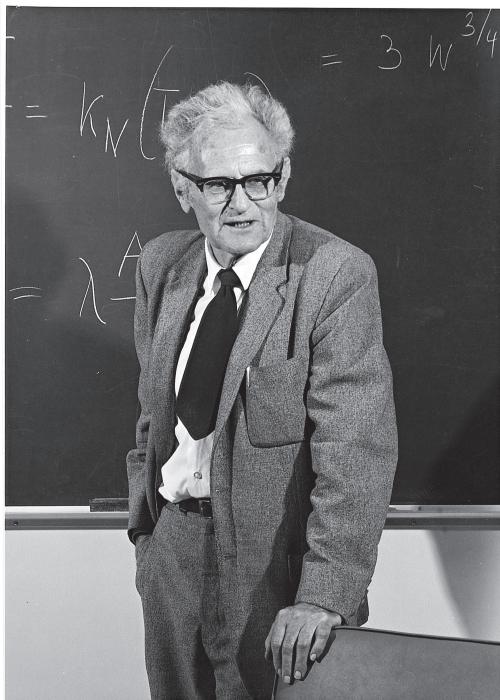


Figure 2.1. Professor Max Kleiber (1893–1976) of the University of California at Davis, Department of Animal Husbandry.

means “other than metric,” or other than linear.) Although size does matter, we don’t live in a world where, e.g., our strength, as measured by the weight we can lift, varies (“scales”) linearly with our own weight. If this were the case, we would be able to lift 50 times our body weight, as some ants are able to do. (Presumably, the *isometric* scaling on Planet Krypton is what endows Superman with his superhuman strength, or so reasoned the creator of this character.) Or, the fact that a cat is 100 times more massive than a mouse would mean that the cat would live 100 times as long. These we know not to be true. If biology does not scale linearly with size, how does it scale?

In 1932, Max Kleiber (Figure 2.1) plotted the logarithm of mass ( $M$ ) in kilograms of various animals and birds on one axis and the logarithm of their basal (resting) metabolic rate ( $Y$ ) (the amount of calories they consume each day) in kcal/day on another axis, and beheld an amazing sight. Figure 2.2 is a recent reproduction of Kleiber’s plot by West and Brown (2004) in *Physics Today*. From dove to hen, from rat to man, from cow to steer, across four orders of magnitude (a factor of  $10^4$ ) difference

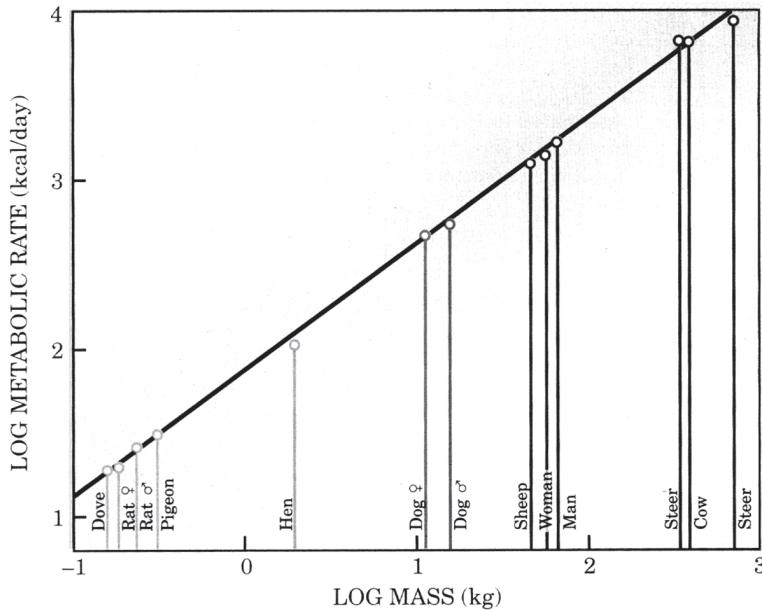


Figure 2.2. The basal metabolic rate of mammals and birds was originally plotted by Max Kleiber in 1932. In this reconstruction, the slope of the best straight-line fit is 0.74, illustrating the scaling of metabolic rate with the  $\frac{3}{4}$  power of mass. The diameters of the circles represent estimated data errors of  $\pm 10\%$ . We use the notation of  $\log$  to denote logarithm of base 10, so that  $\log(10^a) = a$ . Taken from West and Brown, *Physics Today*, September 2004.

in mass, the data lie on a single curve—a straight line in fact! The slope of this straight line is  $b \sim 0.74$ , or about  $\frac{3}{4}$ .

Mathematically, Kleiber showed that

$$\log Y = b \log M + C.$$

This implies that

$$Y = Y_0 M^b \quad (2.1)$$

for some constant  $Y_0 = 10^C$  (which may be different for different taxonomic groups of organisms), implying that there is a power-law dependence on mass  $M$ . Because the exponent  $b$  is less than 1, larger animals tend to have slower metabolic rates. A mouse must eat half its body weight each day so as not to starve to death, while the human's daily consumption is only 2% of his or her body weight. The difficult

part is to come up with a theory for why the power (exponent) is a simple multiple of a quarter, i.e.,  $b \sim \frac{3}{4}$ . This “law of quarters” went unexplained for 50 years.

In 1883, a German physiologist, Max Rubner, actually proposed a scaling law and explained it thus: if an animal is  $L$  times taller than another, then its surface area should be  $L^2$  greater and its mass  $L^3$  greater, since mass is proportional to volume. Its metabolic rate, then, which depends on the amount of heat it sheds, should vary according to its surface area,  $L^2$ , which is proportional to  $M^{2/3}$ . It at least explains why the power  $b$  should be less than 1; otherwise a mouse scaled to the size of a cow linearly would burn itself to death by the amount of heat its large body would now generate but that its surface would be unable to dissipate. Unfortunately this law of thirds did not hold up, and Kleiber appeared to have proved it wrong, although this matter is still being debated. In his book on biological scaling, Schmidt-Nielsen (1984) concluded that “the slope of the metabolic regression line for mammals is 0.75 or very close to it, and definitely not 0.67.”

There are many such scaling laws with the mysterious quarter powers. The lifespan of mammals was found to scale with their mass as  $M^{1/4}$ . Larger animals tend to live longer, but not as long as they would if the lifespan scaled linearly with their size. Larger mammals also have slower heart rates. Their heart rates scale with their mass as  $M^{-1/4}$ . So, the product of lifespan and heart rate, which is the total number of heartbeats in an animal’s lifetime, is independent of the size of the animal. It is about 1.5 billion heartbeats. That is how many heartbeats each of us has in total. It seems that when we use it up, we die. A mouse just uses it up faster than a cat. A cat, which is 100 times more massive than a mouse, then lives  $(100)^{1/4} \sim 3$  times as long as a mouse.

### 2.3 A Model of Branching Vascular Networks

Professor James H. Brown is a field animal biologist at the University of New Mexico. He wrote on this subject and lectured in his class about the quarter-power scaling law for animals of various sizes. His graduate student Brian Enquist wondered if the same law applied to plants. For his PhD dissertation, Enquist discovered the equivalent of Kleiber’s law for plants. So, there seems to be a universal scaling law, whether or not the organism has a heartbeat.

Brown and Enquist suspected that the vascular networks that transport nutrients in plants and oxygen in animals share some commonality in the way they scale with size. They realized, however, that they needed more mathematical power in order to solve this problem. In 1995, Brown and Enquist got together with Geoffrey West, a physicist

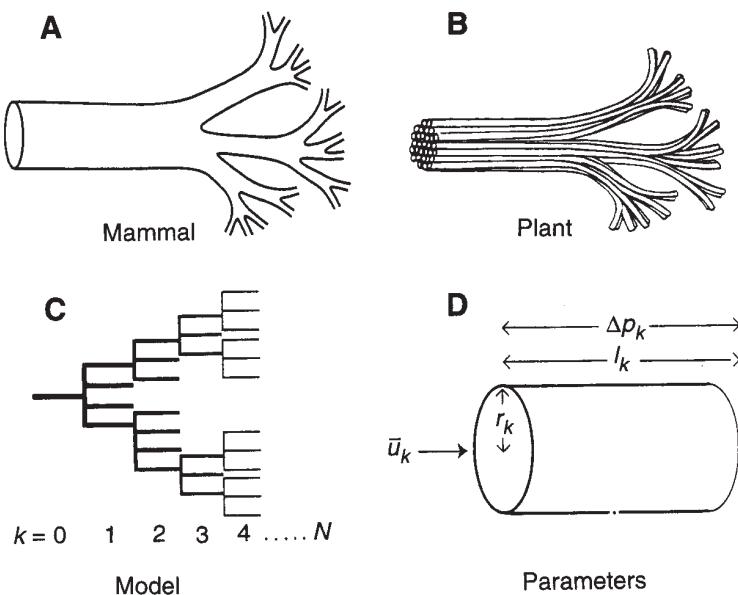


Figure 2.3. (A) Mammalian circulatory and respiratory systems composed of branching tubes; (B) plant vessel-bundle vascular system composed of diverging vessel elements; (C) a schematic representation of such networks, where  $k$  specifies the order of the level, beginning with the aorta ( $k = 0$ ) and ending with the capillary ( $k = N$ ); and (D) parameters of a typical tube at the  $k$ th level. Taken from West, Brown, and Enquist (1997).

and mathematician at the Santa Fe Institute in New Mexico. Their groundbreaking paper on “A General Model for the Origin of Allometric Scaling Laws in Biology” was published in the leading journal *Science* in 1997 (West, Brown, and Enquist, 1997). The mathematics used in that paper is more advanced than the prerequisites for this book. Here we present a simplified derivation. It should be kept in mind that this model still has some loose ends and is still controversial.

The trio focused first on what all animals and plants have in common: a branching vascular network for distributing nutrients or oxygen. In Figure 2.3 we show one such network for a mammal in (A) and for a plant in (B). In humans, our vascular network transports blood with oxygen in a branching network starting with the aorta ( $k = 0$ ) and ending in the capillary ( $k = N$ ) in  $N \sim 22$  levels, with a branching ratio of about  $n \sim 3$  (the number of daughter branches arising from one parent branch). At the capillary level, the nutrients are delivered to the cell tissues serviced by each capillary and wastes are collected.

The authors proposed that the properties of the smallest units in such a network, the capillaries, are the same for animals of different sizes. Thus the capillaries of an elephant have the same width as those of a mouse, under this assumption. Animals of smaller sizes than humans have fewer branching levels ( $N$  smaller than 22), while larger-sized animals have even more branching levels. As a result of natural selection, a similar design is used at different scales. This is called *self-similarity*. A self-similar picture is one in which, if you zoom in on a smaller portion of the picture and magnify it, it will look just like the bigger picture. West, Brown, and Enquist called this property “self-similarity fractal” and they used the phrases “self-similar” and “fractal” interchangeably, although strictly speaking what they had in mind was only “fractal-like.” To be truly “fractal” the self-similarity should be extendable to smaller and smaller scales, and not terminate at the capillary level.

Specifically, the self-similarity assumption says that the ratio of the widths of the vessel from one level to the other:

$$\beta_k \equiv r_{k+1}/r_k = \beta,$$

the ratio of their lengths:

$$\gamma_k \equiv l_{k+1}/l_k = \gamma,$$

and the ratio of their number of branches from one level to the other:

$$n_k = n,$$

are all independent of  $k$ . These “fractal properties,” as the authors called them, are in fact related. As one vessel branches into  $n$  daughter vessels, the cross-sectional area should be preserved to allow a smooth fluid flow. This is called the *area-preserving principle*. Therefore:

$$\pi r_k^2 = n \pi r_{k+1}^2. \quad (2.2)$$

This then implies that

$$r_{k+1}^2/r_k^2 = n^{-1},$$

or

$\beta = n^{-1/2}.$

(2.3)

An additional assumption is called the *volume-filling principle* and is somewhat more difficult to understand because of the ambiguous way it was stated by West et al. This led to some controversy in the scientific literature later (see exercise 2). Essentially, this principle ensures that all cells in a body are serviced by the capillaries. A single capillary of length  $l_N$  and radius  $r_N$  services the surrounding cells of volume  $\pi\rho^2l_N$ , where  $\rho \gg r_N$  is the radius of this cell-cylindrical volume surrounding a single capillary vessel and receiving nutrients from it. It is assumed that

$$\rho \propto l_N.$$

Thus the total volume  $V$  serviced by all the capillaries (whose number is  $N_N$  in their clumsy notation) is

$$V = \pi\rho^2l_N N_N \approx Cl_N^3 N_N, \quad (2.4)$$

where  $C$  is a constant of proportionality. This total volume is the volume of the body of the animal. This same body is also serviced by the vessels in the next level up from the capillary. So it then follows that

$$Cl_k^3 N_k \approx Cl_{k+1}^3 N_{k+1}$$

(as long as  $k$  is large). From it one obtains

$$l_{k+1}^3/l_k^3 \approx N_k/N_{k+1} = n^{-1}.$$

It then follows that

$$\boxed{\gamma \approx n^{-1/3}.} \quad (2.5)$$

Now consider the total volume of blood contained in all the vessels:

$$V_{\text{blood}} = N_0 V_0 + N_1 V_1 + \cdots + N_N V_N,$$

where  $V_k = \pi r_k^2 l_k$  is the volume contained in each blood vessel at level  $k$ . The number of blood vessels at level  $k$  is  $N_k = n^k$ . Using self-similarity, this sum can be written in the form

$$\begin{aligned} V_{\text{blood}} &= \pi[r_0^2 l_0 + n r_1^2 l_1 + \cdots + n^N r_N^2 l_N] \\ &= V_0[1 + (n\beta^2\gamma) + \cdots + (n\beta^2\gamma)^N]. \end{aligned}$$

This is a *geometric series* in the form of

$$S = 1 + r + r^2 + r^3 + \cdots + r^N,$$

where each succeeding term in the sum is a constant factor  $r$  of the term preceding it. If you have forgotten how to sum a geometric series, here is the trick. Multiply  $S$  by  $r$ :

$$rS = r + r^2 + \cdots + r^N + r^{N+1}.$$

Subtract it from  $S$  and observe that all the middle terms cancel:

$$S - rS = 1 - r^{N+1}.$$

Thus the sum of the geometric series is obtained as

$$S = \frac{1 - r^{N+1}}{1 - r}.$$

Using this formula, we can now sum up the blood volumes in the vascular network:

$$V_{\text{blood}} = V_0 \frac{1 - (n\beta^2\gamma)^{N+1}}{1 - (n\beta^2\gamma)}.$$

Because  $(n\beta^2\gamma) = n^{-1/3}$  is less than 1, raised to a large power  $(N + 1)$  it will be reduced to a very small number. Therefore, approximately:

$$V_{\text{blood}} \cong V_0 / [1 - (n\beta^2\gamma)]. \quad (2.6)$$

Since the volume of one capillary vessel is

$$V_N = \pi r_N^2 l_N = V_0 (\beta^2 \gamma)^N,$$

we can rewrite (2.6) as

$$V_{\text{blood}} \cong V_N \frac{(\beta^2 \gamma)^{-N}}{1 - (n\beta^2\gamma)}. \quad (2.7)$$

## 2.4 Predictions of the Model

The total volume of blood in an animal scales with its size and hence mass  $M$ , while the volume of one capillary is independent of size by our prior assumption that the smallest units in our network, the capillaries, are invariant. Equation (2.7) implies that

$$(\beta^2 \gamma)^{-N} \propto M$$

or

$$(\beta^2 \gamma)^{-N} = M/M_0$$

for some constant of proportionality  $1/M_0$ . Taking the log on both sides, we find that, using Eqs. (2.3) and (2.5):

$$N = -\frac{\log(M/M_0)}{\log(\beta^2 \gamma)} = (3/4) \frac{\log(M/M_0)}{\log(n)}. \quad (2.8)$$

The fact that the total number of branches in an animal is proportional only to the logarithm of its mass implies that this fractal-like design is quite efficient for the larger animals. A whale is  $10^7$  times more massive than a mouse but has only 7 times more branchings from aorta to capillary.

The total number of capillaries is given by

$$N_N = n^N = (M/M_0)^{3/4}. \quad (2.9)$$

The metabolic rate  $Y$  is proportional to the rate of flow of nutrients through the vascular network,  $Q_0$ , where if we let  $Q_k$  be the fluid flow rate through each  $k$ th-level vessel, then  $Q_0$ , being that through the aorta, would be the total fluid flow rate. The total flow rate can also be calculated at each level, and it should be  $Q_k$  times the total number of such vessels,  $N_k$ . Thus conservation of fluid flow through each level implies

$$Q_0 = N_k Q_k = N_k \pi r_k^2 u_k = N_N \pi r_N^2 u_N, \quad (2.10)$$

where  $u_k$  is the mean velocity of the fluid through the  $k$ th-level vascular vessel. Given our scaling for  $N_k$  and  $r_k$ , (2.10) implies that the fluid velocity is constant through the vascular network all the way to the capillary, and since the latter is invariant with respect to size, the fluid velocity with which the nutrient is delivered is also invariant. Equations (2.10) and (2.9) imply

$$Y \propto Q_0 \propto N_N = (M/M_0)^{3/4}. \quad (2.11)$$

We have thus derived Kleiber's  $\frac{3}{4}$  power law (2.1) for the metabolic rate.

Equation (2.10) also implies that the radius of the aorta vessel,  $r_0$ , should scale with mass as  $M^{3/8}$  (exercise 1), very close to the observed scaling of 0.36. The length of the aorta,  $l_0$ , should scale as  $M^{1/4}$ .

Another consequence of this scaling with size is explored in West, Brown, and Enquist (2001). They found that if time is normalized by the quarter power of the mass of the mature animal, its growth in time for different animals, be they mammals, birds, fish, or crustaceans, is described by the same universal curve, as seen in Figure 2.4. A derivation of this “universal law of growth” is left to exercise 4. In particular, it is shown that the interval between heartbeats should scale as  $M^{1/4}$ .

For plants, the rate of resource used by each plant (its metabolic rate  $Y$ ) is dependent on its size as  $M^{3/4}$ . Therefore, in an environment with a fixed supply of resources, the maximum number  $N_{\max}$  of plants of average size  $M$  must scale as  $M^{-3/4}$ . Plant ecologists commonly use  $M$  as the dependent variable, and thus

$$M \propto N_{\max}^{-4/3}.$$

This is the so-called  $-\frac{4}{3}$  law in plant ecology.

## 2.5 Complications and Modifications

The *area-preserving principle* of (2.2) is more valid for plant vascular systems than it is for animals with a beating heart. The terminal units of a plant's vascular network are the leaves. It is the transpiration at the leaves that creates the vapor pressure that sucks nutrients through its fiber vessel from the roots, up the trunk and branches, and eventually to the leaves. This “pumping” process is helped along by the osmotic pressure. The vessel tubes are of approximately the same diameter

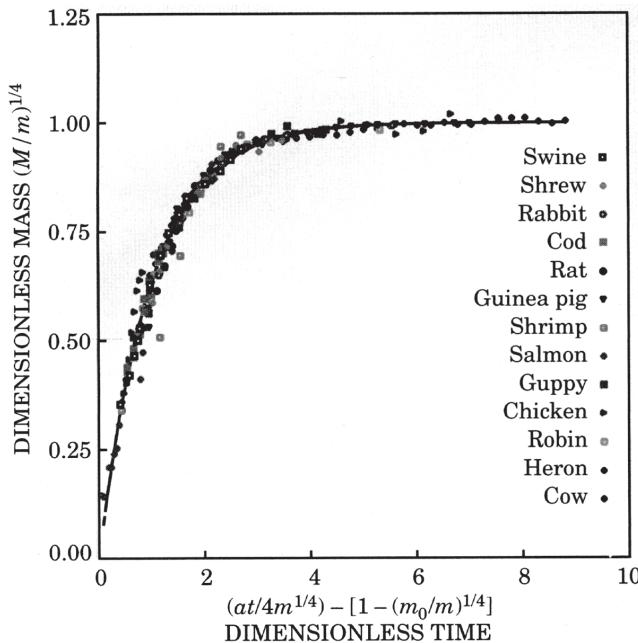


Figure 2.4. The universality of growth is illustrated by plotting a dimensionless mass variable against a dimensionless time variable. Data for mammals, birds, fish, and crustacea all lie on a single universal curve. The quantity  $M$  is the mass of the organism at age  $t$ ,  $m_0$  its birth mass,  $m$  its mature mass, and  $a$  a parameter determined by theory in terms of basic cellular properties that can be measured independently of growth data. Taken from West and Brown, *Physics Today*, September 2004.

throughout this vascular network, but they are bundled together to form the vessel bundles, which then split up into smaller and smaller bundles from the trunk (the “aorta”) to the branches (the “arteries”). Cross-sectional area is preserved in this arrangement (see Figure 2.3B). A consequence of the area-preserving principle is the constancy of the fluid velocity within the vessels (from Eq. (2.10)). For the mammalian vascular network, two complications arise. First, the flow, pumped by a beating heart, is pulsatile. Second, as the branching proceeds from large to small tubes, from aorta to large arteries to smaller arteries, viscosity becomes important, and the flow slows down, which is inconsistent with the area-preserving principle. West, Brown, and Enquist (1997) incorporated these complications and showed that the scaling exponent is dominated by the larger arteries, which are less affected by viscosity

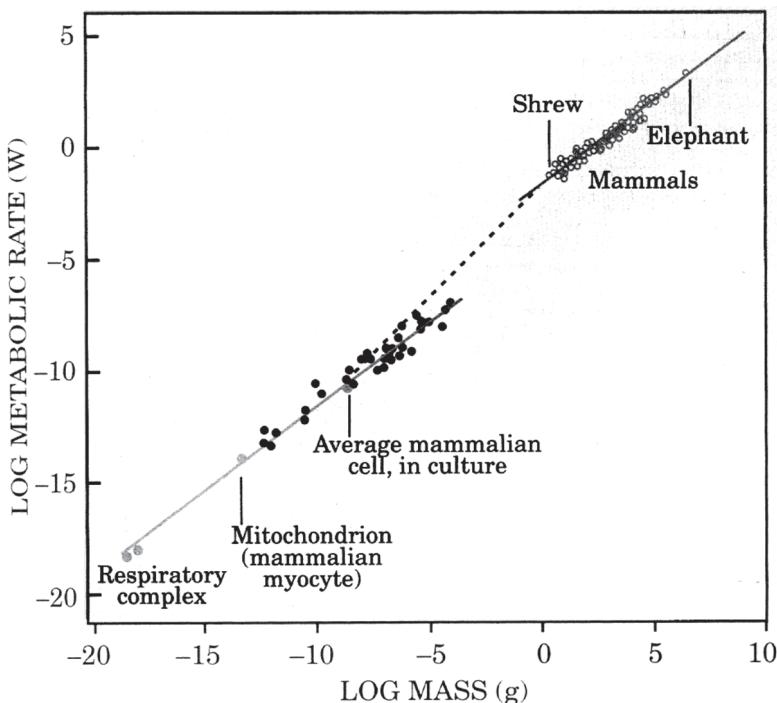


Figure 2.5. The  $\frac{3}{4}$  power law for the metabolic rate as a function of mass is observed over 27 orders of magnitude. The masses covered in this plot range from those of individual mammals, to unicellular organisms, to uncoupled mammalian cells, mitochondria, and terminal oxidase molecules of the respiratory complex. The solid lines indicate  $\frac{3}{4}$  power scaling. The dashed line is a linear extrapolation that extends to masses below that of the shrew, the lightest mammal. Taken from West and Brown, *Physics Today*, September 2004.

and therefore still satisfy the area-preserving principle. Therefore, the previously obtained  $\frac{3}{4}$  power law still holds approximately. For small mammals, such as a 3-gram shrew, the viscous flow starts to dominate just beyond the aorta. There is some indication that smaller mammals deviate from the  $\frac{3}{4}$  scaling, possibly for this reason.

## 2.6 The Fourth Fractal Dimension of Life

Amazingly, the  $\frac{3}{4}$  power law holds not only for plants and animals that have a branching vascular network to carry nutrients but also for single-cell (called unicellular) organisms as well (see Figure 2.5). Thus the  $\frac{3}{4}$  power scaling law appears to span 27 orders of magnitude in mass.

What is the basis for this “universal scaling law of life”? Certainly the model of branching networks needs to be generalized. The *Science* paper by West, Brown, and Enquist (1999) attempts such a generalization by using concepts from fractal geometry and general but less mechanism-specific statements such as, “Natural selection has tended to maximize both metabolic capacity, by maximizing the scaling of exchange surface areas, and internal efficiency, by minimizing the scaling of transport distances and times. These design principles are independent of detailed dynamics and explicit models and should apply to virtually all organisms.” And, “Fractal-like networks effectively endow life with an additional fourth spatial dimension.” The basic idea behind this proposal is that the surface that matters is not the smooth skin enclosing the exterior, as in Max Rubner’s 1883 argument, but the internal surface across which exchange of nutrients takes place. Evolution has selected organisms that maximize the internal surface inside a compact external volume. It leads to the internal surfaces filling up the volume, endowing it with apparently more of the dimensions of a volume (3) than the dimensions of the exterior surface (2). This idea has not yet been commonly accepted, and we await further development in the field. Not wanting to draw too much attention to it, I have left the derivation of this result to exercise 5, where the idea of a fractal dimension is also discussed.

## 2.7 Zipf's Law of Human Language, of the Size of Cities, and Email

Zipf's law is named after the Harvard linguistics professor George Kingley Zipf (1902–1950). Zipf was interested in uncovering the fundamental law of human language. He was independently wealthy and used his own money to hire a roomful of human “computers,” who would count the number of times an English word, such as “and” or “the,” occurs in a given text. He found a power law relating the frequency of occurrence of a word and the rank of that word. The power law is of the form of Eq. (2.1), with the exponent equal to approximately  $-1$ . That is, let  $k$  be the rank of a word. For example, in most English texts, the word “the” is the most often used and therefore is given the rank  $k = 1$ . The second most used word is probably “of,” and it is given the rank  $k = 2$ , etc. Let  $f_k$  be the frequency of occurrence of each word of rank  $k$ , i.e., how many times this word appears in a given text. Zipf's law then says that

$$f_k \propto k^{-b} \quad (2.12)$$

with  $b \sim 1$ .

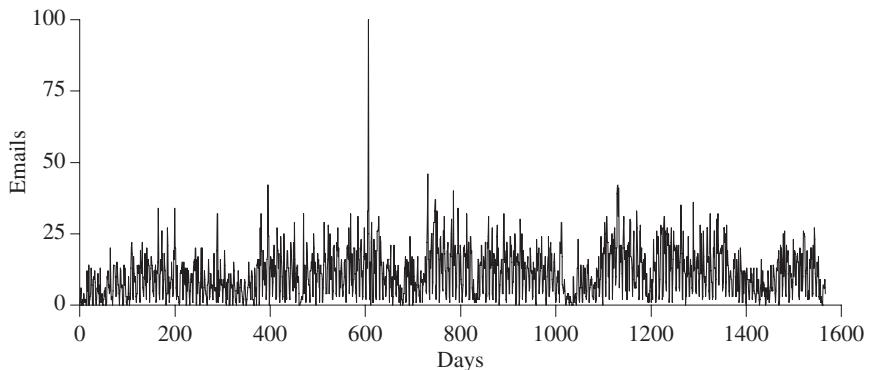


Figure 2.6. The number of emails Professor Mark Kot received each day.

If Zipf's law holds for these  $k$ 's then "of" will appear half as often as "the." This is approximately true in the *Brown Corpus of Standard American English*, a compilation of American English from various sources of one million words, where "the" occurs 7% of the time and "of" 3.5% of the time. For different texts the rank of a particular word may be different. Analyzing the different Zipf curves for texts written by different people may help identify authorship and uncover evidence of plagiarism.

Zipf also analyzed the sizes of cities, as measured by their populations, plotted them against the ranks of the cities, and showed that they also satisfy a power law with  $b \sim 1$ . In the United States, the largest city is New York City, with rank  $k = 1$ , followed by Los Angeles with  $k = 2$  and Chicago with  $k = 3$ , and so on. Plotting the size against rank will also yield a Zipf curve of the form of (2.12). This law appears to apply to most developed countries, but not so well to countries with unique political, economic, or cultural constraints on the movement of its citizens.

My colleague Professor Mark Kot noticed that the number of emails he has received over the past few years appears to follow Zipf's law as well. Figure 2.6 shows the number of emails he received each day. As is typical of most academics, Kot's email number seems to follow an academic-year cycle, with a pause every seven days (on Sunday). When plotted in a log-log plot against the rank of the sender, it approximates a straight line with slope  $-1$  (Figure 2.7). At the low end of the curve, there was a large number of senders who sent only one email, and a smaller number who sent two pieces. The sender ranked number 1, represented by a dot on the high end of the curve, was me. Professor Kot remarked that although I sent the most email to him, I did not send as much as I should have to satisfy Zipf's law. Overall, however, Zipf's law appears to be approximately true here, and we would expect that with

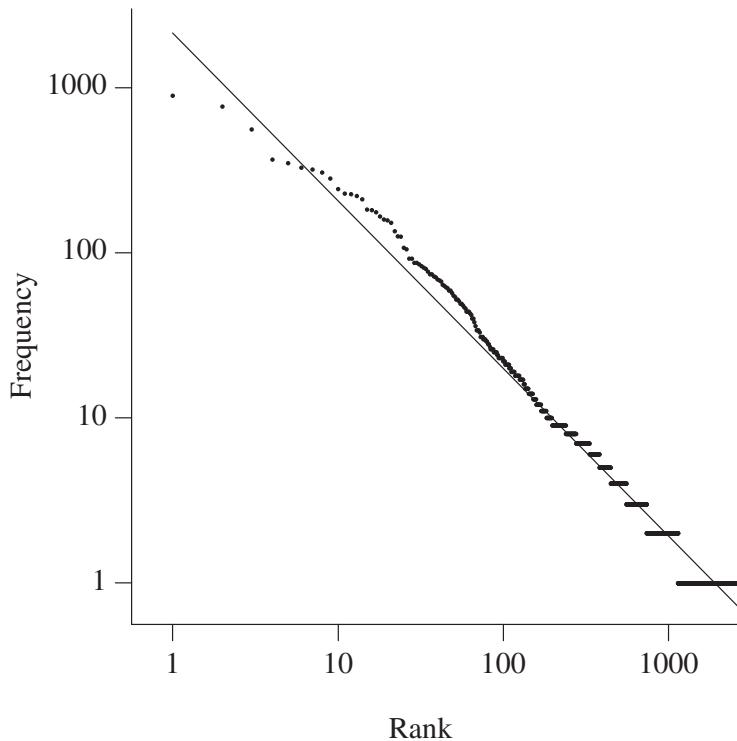


Figure 2.7. Log-log plot of the number of emails against the rank of the sender.

more data, the resulting curve should look better because it may better average out a few individuals' particular behaviors.

Kot, together with his biologist wife, Celeste Berg, and his former PhD student Emily Silverman, looked at the postings in several open biology newsgroups. Submissions (postings) to two such newsgroups are shown in Figure 2.8 as a function of rank of the sender. An almost straight line with a negative slope  $b$  is found, implying a power law in the form of (2.12), though  $b$  can be different from 1. Their finding, together with an explanation of the observed results, is published in Kot, Berg, and Silverman (2003).

There are various proposed models of varying degrees of complexity that try to explain these scaling laws. For example, Marsili and Zhang (1998) showed that Zipf's law for city distribution can be explained by assuming that citizens interact with each other in a pairwise fashion in choosing a city in which to reside. Kot, Berg, and Silverman (2003) proposed a stochastic (probabilistic) model for the scaling law of newsgroups and explained it by computer simulation. Underlying both models is the phenomenon that "the rich get richer."

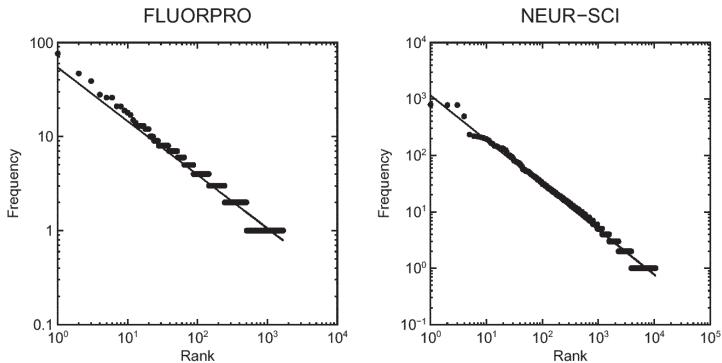


Figure 2.8. Log-log plot of the frequency of submissions as a function of the rank of the sender for two newsgroups, FLUORPRO and NEUR-SCI. Taken from Kot, Berg, and Silverman (2003).

On the other hand, Zipf's original law on word frequency, which he thought revealed the deep law of human language, turns out to be rather "shallow." Wentian Li of the Santa Fe Institute showed in 1992 that a monkey randomly hitting keys on a typewriter with  $M$  letters and 1 blank space generates a random text whose word frequency vs. rank obeys Eq. (2.12), with the exponent  $b$  given by  $\ln(M + 1)/\ln(M) \sim 1$ , with weak dependence on  $M$  as long as  $M > 1$  ( $b = 1.01158$  for  $M = 26$  in English). The derivation is left to exercise 6.

## 2.8 The World Wide Web and the Actor's Network

In 1998, physicist Albert-László Barabási and his colleagues Hawoong Jeong and Reka Albert at the University of Notre Dame embarked on a project to map the World Wide Web, using a virtual robot to hop from one web page to another and collect the links to and from each web page. Counting how many web pages have  $k$  links, they were surprised to find that while more than 80% of them have fewer than four links, 0.01% of them have more than a thousand. Some even have millions. The probability (or the fraction) that any node is connected to  $k$  other nodes was found to be proportional to  $k^{-\gamma}$ , a power law. When plotted on a log-log scale, it follows a straight line with a negative slope of about  $\gamma = 2.1$  (see Figure 2.9B). Similar power laws were found in large social networks, such as the network of actors, where each actor represents a node and two actors are linked if they were cast in the same movie together. The negative slope in this case is  $\gamma = 2.3$  (see Figure 2.9A). For some unexplained reason, the slope seems to lie between 2 and 3 for large networks. The smaller network of power grids in the United States

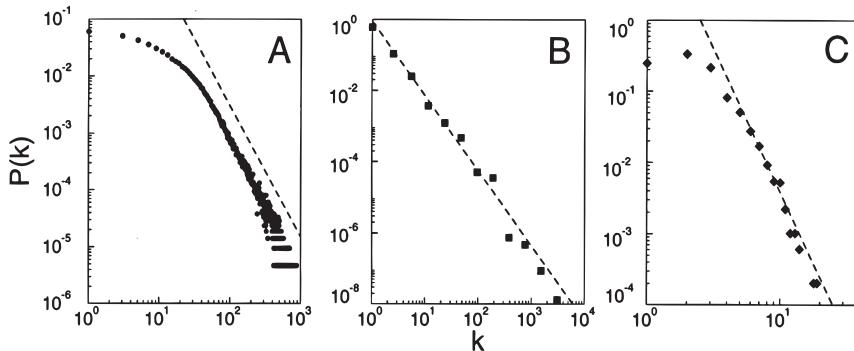


Figure 2.9. The distribution function of connectivities for various large networks. (A) Actor collaboration graph with  $N = 212,250$  vertices, and average connectivity  $\langle k \rangle = 28.78$ . (B) World Wide Web,  $N = 325,729$ ,  $\langle k \rangle = 5.46$ . (C) Power-grid data,  $N = 4,941$ ,  $\langle k \rangle = 2.67$ . The dashed lines have slopes (A)  $\gamma_{\text{actor}} = 2.3$ , (B)  $\gamma_{\text{www}} = 2.1$ , and (C)  $\gamma_{\text{power}} = 4$ . From Barabási and Albert (1999).

(see Figure 2.9C) follows a power law with negative slope of  $\gamma = 4$ . These results were published in *Science* in 1999. That paper, by Barabási and Albert (1999), became highly influential.

The authors were surprised by their finding because they thought people would have different reasons for deciding which sites to link their own web pages to, depending perhaps on the subject matter of their page. Given the diversity of subjects and interests and the vast number of sites one could link to, the situation should appear fairly random. Yet the authors' finding contrasts with the prediction of the standard theory of random networks, which proved that given a number of nodes and random connections between them, the resulting system will have a deeply democratic feature: most nodes will have approximately the same number of links. In the case of power laws, there is no typical number of links, and any number of links is possible. The authors called this type of network "scale free" since it does not have a typical number of links.

There are two important differences between a real network such as the World Wide Web and the textbook case of random networks. The real network is growing, by the continual addition of nodes, and the new nodes preferentially attach to sites that are already well connected. So again the rich get richer, resulting in a few sites (such as Google and Yahoo) with a huge number of links. Early versions of the model concentrated on the development of "scale free" networks, i.e., on explaining why the dependence on the number of links is a power law. Getting the exact value of the exponent was not a high priority, and

the value of the exponent predicted, 3, was deemed good enough. This is in contrast to the case of West, Brown, and Enquist's model of allometric scaling in biological organisms, where the central problem is predicting the exact exponent. Later versions of the network model included more complicated processes of attachment to nodes to get a more specific exponent. We will give here a simplified derivation of the power law, following the earlier model of Price (1976), which preceded that of Barabási and Albert by two decades, although the Price model was about paper citations in journals. A good review can be found in Newman (2003).

## 2.9 Mathematical Modeling of Citation Network and the Web

Derek de Solla Price (1922–1983) was a physicist who later turned into a historian of science. In 1965 he published in *Science* a paper on a pattern of bibliographic citations in science journal papers (Price 1965). It is probably the first example of a network exhibiting a scale-free power-law behavior. An explanation of this power law was later given by Price (1976). A citation network is very similar to a network of web pages. Price assumed that the rate at which a paper gets new citations by other papers is proportional to the number of citations this paper already received. The more highly cited a paper is already, the more likely it is that it will be cited by the next new paper; again the rich get richer. Price assumes that each paper can cite on average  $m$  other papers, with  $m$  fixed. This is approximately true because most journals have limits on the number of references a paper can have. In journals that have no limits, in practice the number of cited references in an average paper does not vary that much. Since each paper on average cites  $m$  other papers, the average number of citations a paper receives within this network is also  $m$ . Let  $N_k$  be the number of papers (nodes) in the network that have been cited  $k$  times, and let  $n$  be the total number of papers in the whole network. We write  $p_k = N_k/n$  as the fraction of the papers in the network having  $k$  citations. As a new paper is published, the network size increases by 1, from  $n$  to  $n + 1$ . We can use  $n$  as a timelike quantity, although it increases by discrete integers only. So the change in  $N_k$  with the arrival of this new paper is given by

$$N_k(n+1) - N_k(n) \equiv (n+1)p_k(n+1) - np_k(n).$$

This change is given by the number of papers that had originally been cited only  $k - 1$  times and are now cited by this new paper, thus becoming papers with  $k$  citations. It is decreased by the number of papers that already had  $k$  citations and that this new citation moves them into the

category of papers with  $k + 1$  citations, therefore no longer belonging to the category of papers with  $k$  citations. The probability that a new paper will cite any paper that already has  $k$  citations is proportional to  $k$ , which is a central assumption of Price, as mentioned above. This assumption, however, has a slight technical difficulty: the new paper has no one citing it yet, and so  $k = 0$ . This assumption then ensures that no future papers will ever cite it. Price got around this difficulty by saying that one should count the citation in such a way that a new paper should come with one citation, that by itself. Then the probability that a new paper will cite any paper that has  $k$  citations is proportional to  $(k + 1)$ , and in fact equal to

$$\frac{(k+1)p_k}{\sum_k (k+1)p_k} = \frac{(k+1)p_k}{(m+1)},$$

since  $\sum p_k = 1$  and  $\sum kp_k = m$ , the mean number of citations in the network. Furthermore, since the mean number of new citations by each new paper is  $m$ , the number of new citations of existing papers already with  $k$  citations is  $m$  times the above probability, and thus  $(k+1)p_k m / (m+1)$ . Therefore,

$$N_k(n+1) - N_k(n) = kp_{k-1}(n) \frac{m}{m+1} - (k+1)p_k(n) \frac{m}{m+1}. \quad (2.13)$$

For large networks ( $n$  large), the fraction  $p_k(n+1) \sim p_k(n)$ . Thus  $N_k(n+1) - N_k(n) \equiv (n+1)p_k(n+1) - np_k(n) \sim p_k(n)$ . Equation (2.13) becomes

$$p_k = \frac{m[kp_{k-1} - (k+1)p_k]}{(m+1)}.$$

This is the same as

$$p_k = \frac{kp_{k-1}}{(k+2+1/m)}. \quad (2.14)$$

This is a first-order difference equation, except that the coefficient is dependent on the variable (the index  $k$ ). Nevertheless, one can simply iterate. Starting with  $k = 1$ , we get from Eq. (2.14)  $p_1 = p_0/(3+1/m)$ . With  $k = 2$ , it is  $p_2 = 2p_1/(4+1/m) = 2p_0/(4+1/m)(3+1/m)$ , and so on. It then yields the following general solution:

$$p_k = \frac{k(k-1)(k-2)\cdots(2)(1)}{(k+2+1/m)(k+1+1/m)\cdots(3+1/m)} p_0. \quad (2.15)$$

It turns out that this unwieldy expression, when plotted, looks like a power law for large  $k$ :

$$p_k \propto k^{-(2+1/m)}. \quad (2.16)$$

There are two ways to demonstrate this. The long fraction in (2.15) is not just any fraction; it is actually a special function, Legendre's beta function, that mathematicians have studied. It is known that for large  $k$  it has the asymptotic behavior of (2.16). Another way to show this is to observe that, for large  $k$ , the difference equation (2.14) can be written approximately as

$$\frac{p_k}{p_{k-1}} \approx \left(1 - \frac{(2+1/m)}{k}\right). \quad (2.17)$$

This is because, using the binomial expansion of  $1/(1+x) \approx 1-x$ ,

$$\frac{1}{(k+2+1/m)} = \frac{1}{k\left(1 + \frac{(2+1/m)}{k}\right)} \approx \frac{1}{k} \left(1 - \frac{(2+1/m)}{k}\right).$$

To solve Eq. (2.17), we use the trial solution

$$p_k = k^\lambda,$$

where the exponent  $\lambda$  is unknown and to be determined in the solution process. When substituted into (2.17), the trial solution yields

$$\frac{k^\lambda}{(k-1)^\lambda} \approx \left(1 + \frac{\lambda}{k}\right) = \left(1 - \frac{(2+1/m)}{k}\right),$$

from which we obtain the exponent to the power law:

$$\lambda = -(2+1/m).$$

Equation (2.16) states that the fraction of papers with  $k$  citations by other papers varies with  $k$  in a power law of the form of  $\sim k^{-\gamma}$ , with a negative exponent of  $\gamma = 2 + 1/m$ , slightly larger than 2. Price thought that this was consistent with what he found in the *Science Citation Index*.

We now return to the model of Barabási and Albert for the World Wide Web links, which is actually simpler than that of Price but is just a little more difficult to understand. In the simple model used by Barabási and Albert (1999), it is assumed that each new addition to the

web initially contains  $m$  links, which are directed both ways. That is, it is assumed that this new site links to  $m$  other sites and that these other sites also have links to it. This is unrealistic but gets around the problem facing Price of the new additions having zero probability of being linked by other nodes. Because of these two-way links, the mean number of links in the network is now  $2m$  instead of  $m$  in Price's model. With this simple modification the problem is solvable exactly. This is left as an exercise (exercise 7). It should be pointed out that Barabási and Albert (1999) were mainly interested in showing that the web links are "scale free," i.e., satisfy a power law. They were not too concerned with the exact magnitude of the exponent. Their result of an exponent  $\gamma = 3$  was not too realistic. That problem was left to later authors. See Newman (2003).

## 2.10 Exercises

### 1. Aorta radius

Show that the radius of the aorta blood vessel scales with the mass  $M$  of the animal as  $M^{3/8}$ .

### 2. A critique

Two Polish scientists, Kozłowski and Konarzewski (2004), wrote a critique of the West, Brown, and Enquist (1997) paper, pointing out that the model is mathematically incorrect. Their criticism centers on the expression of Eq. (2.4) for an animal's body volume:  $V = \pi \rho^2 l_N N_N \sim C l_N^3 N_N$ . Since it is proportional to the mass of the animal, the authors pointed out that it then follows from Eq. (2.4) and the assumption that the properties of the capillaries are invariant with respect to an animal's size that the total number of capillaries  $N_N$  in an animal of mass  $M$  should scale linearly with  $M$ , and not as a  $\frac{3}{4}$  power of  $M$  as in Eq. (2.9). Therefore, the final result, Kleiber's law of  $\frac{3}{4}$  powers derived in Eq. (2.11), is wrong, according to these authors. In rebuttal, Brown, West, and Enquist (2005) wrote:

WBE clearly state that only the characteristics of the capillaries themselves are assumed to be invariant. Nevertheless, K&K incorrectly interpreted this size-invariance to mean that each capillary must supply a constant volume of tissue.... Predicting the scaling of the "service volume" of tissue supplied by a capillary is an integral part of the WBE theory. WBE proved that the service volume increases with body size, as  $M^{1/4}$ .

Well, in the original paper, West, Brown, and Enquist (1997) assumed that the service volume surrounding each capillary of length  $l_N$  is a sphere of diameter  $l_N$ . This then leaves no room for the service volume to increase with body size. In our derivation, we have inserted a constant  $C$  for the service volume in (2.4) so that this volume is  $Cl_N^3$ .  $C$  is independent of  $k$ , the level of branching for an animal, but it may be different for different animals. Allowing  $C$  to depend on  $M$ , show that  $C$  is proportional to  $M^{1/4}$ , thus resolving the apparent inconsistency.

### 3. Self-similarity assumption

Etienne, Apol, and Olff (2006) gave a more general derivation of the  $\frac{3}{4}$  power law that does not require the assumption of self-similarity. Read this article and rederive the  $\frac{3}{4}$  power scaling for the metabolic rate result under this more general condition.

### 4. Universal law of growth

(Note: Solution of this problem requires a knowledge of how to solve a simple first-order ordinary differential equation. If you are unfamiliar with ordinary differential equations, please read the review in Appendix A first.) In a growing organism, metabolism supplies energy to both maintain existing tissues and create new tissues by cell division. Let  $Y$  be the metabolic rate of an organism,  $Y_c$  the metabolic rate of a single cell,  $N_c(t)$  the total number of cells at time  $t$ ,  $m_c$  the mass of a cell, and  $E_c$  the energy required to create a new cell. The cell properties,  $E_c$ ,  $m_c$ , and  $Y_c$ , are assumed to be constant and invariant with respect to the size of the organism. Thus:

$$Y = Y_c N_c + E_c \frac{dN_c}{dt}.$$

Let  $m$  be the total body mass of the organism at time  $t$ , and  $m = m_c N_c$ . (Note that  $N_c$  is the total number of cells in a body and is proportional to mass  $m$ , while the total number of capillaries  $N_N$  is proportional to  $\frac{3}{4}$  power of  $m$ .) From (2.11), we have  $Y = Y_0(m)^{3/4}$ .

- a. Show that the above equation can be written as

$$\frac{dm}{dt} = am^{3/4} - bm,$$

with  $a = Y_0 m_c / E_c$  and  $b = Y_c / E_c$ .

- b. Let  $m = M$  be the mass of a matured organism, when it stops growing (i.e.,  $dm/dt = 0$ ). Find  $M$ , and show that the above equation can be

rewritten as

$$\frac{dm}{dt} = am^{3/4}[1 - (m/M)^{1/4}].$$

- c. Let  $r = (m/M)^{1/4}$ , and  $R = 1 - r$ . Then the above equation becomes

$$\frac{dR}{dt} = -\left(\frac{a}{4M}\right)R.$$

Solve this simple ordinary differential equation and show that a plot of  $\ln(R(t)/R(0))$  vs. the no-dimensional time  $at/(4M^{1/4})$  should yield a straight line with a slope  $-1$  for any organism regardless of its size.

- d. Based on this scaling for time  $t$ , argue that, for a mammal, the interval between heartbeats should scale with its size as  $M^{1/4}$ .

### 5. Fractal dimensions

The intuitive concept of dimensions of objects can be stated in a mathematical form in the following manner: a volume  $V$ , which depends on lengths  $x, y, z, \dots$ , becomes another volume  $V'$  when the lengths are multiplied by the same factor  $\lambda$ . Thus

$$V' = V(\lambda x, \lambda y, \lambda z, \dots).$$

If, e.g.,  $\lambda = 2$ , then  $V'$  will be  $8 = 2^3$  times  $V$  when its length, width, and height are all increased by a factor of 2. This relationship is expressed in the following self-similar scaling relation:

$$V' = V(\lambda x, \lambda y, \lambda z, \dots) = \lambda^d V(x, y, z, \dots).$$

In this case  $d = 3$ , and that gives the dimension of the volume. Similarly, we can scale a surface area as

$$S' = S(\lambda x, \lambda y, \dots) = \lambda^d S(x, y, \dots)$$

and deduce that the dimension of the surface is  $d = 2$ . When defined in this more general way, the dimension of an object could give some nonintuitive values, including fractional numbers—hence the word fractal. For a visual example, see exercise 5 of Chapter 1.

Now consider the internal surface across which nutrients are exchanged in a biological organism. As an example, consider the lung, which has a compact and smooth outer surface (the chest wall), and an irregular internal surface enclosing each of a huge number of small

air sacs. We are interested in determining the fractal dimension of the internal surface. The area, denoted by  $A$ , depends on the structure of the branching network characterized by length  $l_k$  at each level. So we write

$$A = A(l_0, l_1, l_2, \dots).$$

Now we want to see how this area scales if all lengths are multiplied by the factor  $\lambda$ :

$$A' = A(\lambda l_0, \lambda l_1, \lambda l_2, \dots) = \lambda^d A(l_0, l_1, l_2, \dots),$$

and the dimension  $d$  should be 2 for a surface, as can be argued intuitively.

Now comes the twist. In the West, Brown, and Enquist model, all lengths  $l_k$  scale with size, except the terminal unit, which is denoted here by  $l_0$ . This terminal unit is invariant to the scaling. This is not the typical fractal self-similarity, and it is somewhat unfortunate that we need to introduce the fractal definition of dimensions while talking about a concept that is not fractal. The biologically relevant scaling needed is instead:

$$A'' = A(l_0, \lambda l_1, \lambda l_2, \dots) = \lambda^{2+\epsilon} A(l_0, l_1, l_2, \dots).$$

The exponent  $\epsilon$  is no longer 0 because of this strange scaling, and is as yet undetermined. The above form of the scaling is only a hypothesis, because it has not been shown that such a surface can be scaled this way. Similarly, it is assumed that volume scales as

$$V'' = V(l_0, \lambda l_1, \lambda l_2, \dots) = \lambda^{3+\sigma} V(l_0, l_1, l_2, \dots).$$

The exponent  $\sigma$  is as yet undetermined.

- a. Show that as far as the scaling with respect to  $\lambda$  is concerned, surface scales with volume as

$$A'' \propto V''^{\frac{2+\epsilon}{3+\sigma}}.$$

Argue that the superscript " can be dropped, resulting in

$$A \propto V^{\frac{2+\epsilon}{3+\sigma}} \propto M^{\frac{2+\epsilon}{3+\sigma}}.$$

$\epsilon = \sigma = 0$  yields the  $\frac{2}{3}$  power law of Rubner, while  $\epsilon = \sigma = 1$  would give the  $\frac{3}{4}$  power law of Kleiber.

- b. West, Brown, and Enquist (1999) suggested that organisms have evolved to maximize their internal surface within a given volume to facilitate exchange of nutrients, and they claim that this is equivalent to maximizing the exponent  $b = (2 + \epsilon)/(3 + \sigma)$ . This expression does not have a distinct maximum unless one allows negative numbers for  $\sigma$ . So, clearly, more assumptions are needed.

Use the following arguments to show that the maximum is attained for  $\epsilon = 1$  and  $\sigma = \epsilon$ : a volume cannot have a dimension less than a surface; the smallest value that  $\sigma$  can take is  $\epsilon$ . On the other hand, the dimension of the internal surface, being an area, cannot exceed the dimension of an ordinary volume; the largest value  $\epsilon$  can take is 1. Hence derive Kleiber's  $\frac{3}{4}$  power law.

- c. Summarize, in your own words, the assumptions needed to derive Kleiber's law vs. those needed for Rubner's law. Are you satisfied with these assumptions in this fractal context?

## 6. Zipf's law for random texts

Following Li (1992), we let there be  $M$  letters and one space that our random monkey can type, each with probability  $p = 1/(M + 1)$ . A word is formed when a string of letters of any size is preceded and followed by a space, denoted by  $_$  here. The probability of getting the word  $_a_$  is  $p^3$ . The probability of finding a three-letter word, such as  $_bst_$ , is  $p^5$ , and that of an  $L$ -letter word,  $p^{L+2}$ . So the frequency of occurrence of any word with length  $L$  is

$$f_i = cp^{L+2},$$

where  $c$  is a normalization constant. There are  $M^L$  such words with length  $L$ . The constant  $c$  is to be determined by the convention that the frequency of occurrence of all words is normalized to 1:

$$1 = \sum_{L=1}^{\infty} M^L \frac{c}{(M + 1)^{L+2}}.$$

The frequency of occurrence of all words of length  $L$  is

$$f(L) = M^L f_i(L).$$

- a. Show that

$$f(L) = \frac{M^{L-1}}{(M + 1)^L}.$$

- b. To convert word length  $L$  into rank, we note that the formula for  $f_i$  implies that shorter words (those with smaller values of  $L$ ) should occur more frequently and hence have a lower value for the rank  $k(L)$ . The rank is easy to understand (the most frequently occurring word is of rank 1, the next is rank 2, etc.) but it is difficult to define mathematically. Nevertheless, argue that the rank can be expressed as

$$k(L) = \alpha M^L - \beta,$$

for some constants  $\alpha$  and  $\beta$ . This yields, upon taking the log with base  $M$ :

$$L = \log_M(\alpha^{-1}(k + \beta)).$$

Substituting this expression for  $L$  into the formula in (a), show that  $f$  satisfies the following generalized Zipf's law:

$$f = C(k + \beta)^{-b},$$

where

$$b = \frac{\ln(M + 1)}{\ln(M)}.$$

Show that  $b \sim 1$ .

### **7. A mathematical model for the World Wide Web**

The model of Barabási and Albert considers the situation when a new node attaches to the existing network consisting of  $n$  nodes. This new node has  $m$  undirected links, meaning that it is linked to  $m$  existing nodes in two directions. If the entire web is built up this way, then the mean number of links is  $2m$ . The probability of attachment by this new node to any existing node that already has  $k$  links is

$$\frac{kp_k(n)}{2m},$$

where  $p_k(n)$  is the fraction of nodes in the network of  $n$  nodes that have  $k$  links. So the number of nodes that gain one link by this additional node is  $m$  times this probability.

- a. Derive the counterpart of Eq. (2.13) by arguing that its right-hand side should be

$$\frac{1}{2}(k - 1)p_{k-1}(n) - \frac{1}{2}kp_k(n).$$

- b. For a large network ( $n$  large), show that this reduces to

$$\frac{p_k}{p_{k-1}} = \frac{k-1}{k+2}.$$

- c. Solve the first-order difference equation above by iteration and show that the solution is a constant times

$$p_k \propto \frac{1}{(k+2)(k+1)(k)}.$$

Therefore, for large  $k$  it is a power law with  $\gamma = 3$ .