

The use of accounting data in statistical models

Duarte Trigueiros

Faculdade de Economia, Universidade de Algarve

Resumo

É difícil usar dados contabilísticos em modelos estatísticos pois, abundantemente, os casos atípicos e as distribuições são heteroscedásticas, leptocurticas e fortemente assimétricas. O objetivo deste artigo é fornecer sugestões que facilitem a compreensão do comportamento estatístico desses dados de modo a permitir uma mais fácil modelação e teste de hipóteses.

Palavras-Chave: Dados contabilísticos, modelação estatística

Abstract

It is difficult to use accounting data in statistical models. Outliers abound, influential cases distort estimation while heteroscedasticity, fat tails and asymmetric distributions make *P*-values meaningless. When trying to solve these problems, practitioners find out that accounting data is not only difficult to use, its statistical behaviour is also difficult to predict.

The goal of this note is to facilitate understanding of the statistical characteristics of accounting numbers. Guidelines are offered on how to transform raw data into well-behaved variables, how to deal with non-proportionality, how to accurately model firm size and solve other difficulties.

Keywords: Accounting data, statistical modelling

1 Accounting Numbers are Multiplicative

Accounting data is the numerical information found in annual reports of firms. Reports contain sets of accounts such as the Profit and Loss Account, the Balance Sheet and others. Each item in these sets reports on a specific magnitude. The volume of sales of the year, for instance, is reported in the item named Sales. Magnitudes found in sets of accounts are the raw material for financial analysis and all types of statistical manipulation. But before being of any use, these numbers are usually combined to form ratios. Typically, two numbers from the same report, say Earnings and Net Worth may be chosen as the numerator and denominator of a ratio. Most accounting data used in statistical models is in the form of ratios.

It is impossible to understand the statistical characteristics of ratios without understanding first the characteristics of the numbers they are made of and the way they interact. Therefore, the following lines are devoted to such numbers.

The first and most important fact about numbers found in accounts is that they cannot be described as resulting from the type of random mechanism that leads to Normal variables. Normal variables stem from additive random mechanisms, an observed distribution is the result of adding a large number of other distributions. For instance, the distribution of weight in adult steers stems from adding the probabilities associated with genetic effects, eating habits and other effects. In the limit, any addition of probability distributions, no matter which, leads to the Normal distribution (this is known as the Central Limit theorem).

Obviously, in an additive mechanism, each intervening effect may lead to an increase or to a decrease in the likelihood associated with the resulting event. Adequate eating habits may, for instance, be able to balance a genetic pre-disposition to put on weight, lowering the expected weight. The mechanism generating accounting numbers is different, effects always reinforce each other. This is because such effects are, in this case, every individual transaction contributing to a reported magnitude. Indeed, each transaction contributing to the amount reported as, say, Sales for a given period, is itself a random event. It contributes to the reported number, say in a manner, which could lead to either an increase or decrease in Sales, but by accumulation only.

Accumulations of random events lead to multiplicative, as opposed to additive variables, because the likelihood of realisations is conditional on the occurrence of a chain of previous events, not on any free interaction of influences, some positive, others negative. Such likelihood thus stems from multiplying rather than adding probabilities.

Multiplicative distributions are easy to recognise. They are skewed, exhibiting long tails towards positive values. As a consequence, some of the observations in a sample are likely to exhibit very large magnitudes in comparison with others, thus giving the impression of being outliers.

Contrasting with additive distributions where skewness and kurtosis are independent, in mechanisms of the Multiplicative type both statistics are manifestations of a unique, underlying phenomenon, variability. Therefore, highly volatile variables exhibit markedly skewed and leptokurtic distributions whereas those where variability is small have almost symmetrical, non kurtotic distributions.

Accumulation is just one amongst several processes leading to Multiplicative numbers. Any variable where magnitude x is, on average, proportional to changes dx , will exhibit a Multiplicative distribution. The natural form of the origin of such type of variable is

$$\frac{dx}{x} = \mu dt + \sigma dz \tag{1}$$

where μ is an expected percent change, x is the variable supposed to drive changes in x , dx is a small random disturbance and μdx is supposed to be independent of x . Therefore, in this type of variable, percent changes are additive. Normality, as a limit, is approached by percent change, not by absolute change.

The mechanism depicted in (2) is known as the "Gibrat's Law of Proportionate Effect". It leads to lognormal distributions (distributions where the logarithm of observations is Normal) or to other types of multiplicative distribution. Multiplicative, proportionate, exponential and lognormal are terms variously used to designate the family of skewed distributions with its origin in (1). Archibison and Brown (1957) describe the lognormal distribution.

Sales, Earnings, Assets and other accounting aggregates are since long known, in domains such as Industrial Economics and others, to be multiplicative, i.e. broadly lognormal. In spite of this, until recently the accounting literature has discarded lognormality in accounting data as incompatible with sound reasoning while quoting the influential fallacy introduced by Eisenbeis (1977) or the fact that some ratios are apparently Normal, or even the existence of negative values in some accounts. Excessive skewness and other characteristics of multiplicative variables were interpreted as distortions of normality or, in the case of ratios, as a side effect of non-proportionality (Haines, 1982). In the Accounting research domain, the case for multiplicative mechanisms was made by McLeay (1986) and by Trigueros (1995).

The peculiar characteristics of lognormal variables must be borne in mind in any context involving the manipulation of these variables or their ratios. Lognormality cannot be treated as a simple departure from normality. For coefficients of variation beyond 0.25, skewness and kurtosis are so severe that most observations concentrate in a small region with only a few extreme values spreading out over a wide range. No parametric tool is robust enough to avoid severe distortion when such data is used.

Statistical models are functional forms supposed to reflect inter-relationships amongst effects. Descriptions of inter-relationships amongst effects greatly differ between additive and multiplicative variables. For additive data, distributions are preserved when variables are added or subtracted. This is not the case for multiplicative data where distributions are preserved when variables are multiplied or divided. The addition or subtraction of two Normal variables will be Normal; the product or ratio of two lognormal variables will be lognormal. The simplest additive formulation would be $x = \mu + \epsilon$ where x is explained as effect μ , the expectation, plus a random deviation ϵ . The multiplicative equivalent would be $x = \mu \epsilon$ where x is now explained as the product of a constant magnitude, μ , by a random factor ϵ .

The likelihood that ϵ may stretch beyond two or three standard deviations above or below μ is very small. Therefore, in general, additive variables describe deviations from an average magnitude but they are unable to describe large differences in magnitude. By contrast, in the case of multiplicative formulations, the exponential nature of ϵ leads to likely values of ϵ over a much wider range. The volume of sales of United Biscuits, a firm in the 95th size percentile of its industry, is about 500 times that of firms in the 5th size

percentage. And the observations would never be able to model such large discrepancies. This is why lognormality often denotes size influence whereas normality generally denotes a size-free variable.

2. How to Use Financial Ratios in Statistical Models

The existence of a size influence in magnitudes reported in the accounts of firms led to the use of ratios. Ratios such as Return on Equity, Interest Cover, Debt to Net Worth and many others are widely used by managers, practitioners, and analysts. They control for size so that comparisons may be made.

When accounting numbers are lognormal, then ratios should be lognormal as well. But some ratios show unexpected characteristics, which makes them difficult to use in statistical models. For instance, although most ratios are indeed lognormal, Total Debt to Total Assets, Net Worth to Total Assets and others are apparently Normal. Current Assets to Total Assets is negatively skewed (see, e.g., So, 1987). How is this possible? The reason is straightforward. Accounting identities preclude some ratios from taking on the values a skewed distribution would allow.

This constraining effect is clearly observable when plotting on a log-ratline scale, the two components of a ratio. Figure 1 shows the constraint imposed by Total Assets on the spread of Net Worth.

Figure 1 Bivariate distribution of Total Assets (X-axis) and Net Worth (Y-axis)



Adequate transformations can take into account constraining mechanisms yielding unimodal ratios. For example, any ratio where the numerator cannot be larger than the denominator (i.e., ratios of the form

$$\frac{1}{\sum_i}$$

can be transformed into the corresponding unconstrained ratio

$$\sum_i \frac{1}{x_i}$$

The unconstrained ratio corresponding to the ratio Fixed Assets to Total Assets (FA/TA) is the ratio FA/CA where CA=TA-FA. The information contained in both ratios is the same. The difference between them is just functional.

Table 1 summarises transformations able to bring ratios affected by several types of constraints into parametric behaviour. Ratios where there is no constraint are logarithmic.

Table 1: Transformations adequate to be used in constrained ratios (from McLeay and Frigueiras, 2003).

Case No.	Ratio	Example	Transformation to use	Boundaries of the ratio
1	R = Y/X	Current Ratio	Log R	0; ∞
2	R = (X-Y)/X	Sales Margin	Log (1-R)	-∞; 0
3	R = (Y-X)/X	Change in Capital Employed	Log (1+R)	-1; ∞
4	R = (Y-X)/X	Interest Cover	Log (R-1)	1; ∞
5	R = X/(Y-X)	Liabilities Ratio	Log (1/(R-1))	0; 1
6	R = X/(Y-X)	Leverage Ratio	Log (1/(R+1))	0; ∞

Awareness of the existence of constraining mechanisms and the way they affect ratios removes one major obstacle in understanding and using financial ratios. But not all is explained. A fact that remains unaccounted for is the existence of fat tails (leptokurtosis) in the logarithms of all types of ratios that is, even after appropriate transformations are applied. Such leptokurtosis, though, is not severe and may, in most cases, be ignored.

Before moving on to the following section it is important to recall two facts about ratios. First, the use of ratios requires caution in order to avoid other, well-documented limitations. Ratios, for instance, produce ambiguous results when denominators, not only numerators, may take on negative values; ratios are sensitive to atypical magnitudes; numbers from the Profit and Loss account must be corrected for reported periods different from one year; accounts from the Balance Sheet may be distorted by seasonal effects. The second important fact is that, notwithstanding the above limitations, ratios do not deserve the negative image conveyed by scientific journals and books. In spite of their widespread use, ratios are regarded with scepticism by scholars and are described as some primitive tradition with no scientific support. The major cause of such scepticism is the mentioned diversity in statistical distributions. Another is the fact that some reported

numbers may take on negative values, which is difficult to reconcile with multiplicative behavior.

The following section discusses the distribution of accounting flows or differences (such as Earnings) which may indeed take on negative values.

3. The Distribution of Earnings

The literature on the distribution of ratios seems to consider profits as additive, albeit not necessarily Normal (McLeay, 1986; Tippett, 1990, amongst others). Probably this is because authors focus on profitability ratios, not on profits.

Profitability ratios basically express percent changes in Net Worth. Authors believe that a ratio of the form $\Delta x/x$ should approach normality since lognormal x lead to Normal $\Delta x/x$ (where Δx are small disturbances). However, in order to qualify as a small disturbance, Δx must be really small when compared with x . Now, for most numbers taken from sets of accounts this is simply not the case. Flows such as changes in Net Worth are not small enough to qualify as a disturbance, not least when compared with the respective stock. Annual Sales, in spite of being a flow, often is larger than stocks such as Fixed Assets, e.g., in high turnover industries. Actually, accounting flows are not that different from stocks, being large accumulations. Thus profitability and other ratios are indeed multiplicative.

Numbers that may take on negative values are simply the result of subtracting positive-only accumulations. A magnitude reported as Earnings is obtained by subtracting the different types of costs and expenses from Revenues. The distribution of these variables should therefore stem from subtracting lognormal distributions. The task of analytically determining such distributions is not easy as it requires working out the logarithm of a subtraction. There is however a fact that simplifies analysis. Costs and Revenues are correlated because both are influenced by the same effect, size. When correlation is taken into account it becomes possible to approximate analytically the distribution of flows. It turns out that, for conditions typically found in industries, size distribution is not unique. Rather, it is a juxtaposition of two approximately lognormal density functions, one for positive and the other for negative values, the latter being a mirror-image of the logarithm as depicted in figure 3. Simulation confirms this result.

Figure 2: The density function of Earnings is a juxtaposition of two lognormal distributions.



Ratios formed with such distributions may be markedly two-tailed, giving the impression that they are near symmetry. Fat-tailed distributions such as Student's t or Cauchy's may indeed fit them closely (McLeay, 1986). It should be made clear, however, that the hypothesis of additive distributions leads to unreasonable conclusions. The reporting of immaterial values must be less likely than that of sign-related values. The probability density of Earnings, for instance, must decrease when approaching zero and then increase again after passing through zero into negative values as predicted by Figure 2. This is because losses, as well as profits, must be proportionate to size. Additive distributions would imply that the reporting of immaterial profits or losses is more likely than that of material losses.

Ratios are multiplicative no matter their components' type, stocks or flows. When the ratio is constrained by some accounting identity, an appropriate transformation will bring it back to a standard behaviour. Specifically, profitability ratios will benefit from transformation 3 in Table 1. Indeed, most commonly found situations involving negative flows are solved simply by using log or other transformation. As for ratios where the denominator, not just the numerator, may take on negative values, it seems as though there is no other choice but to consider two populations, one for positive and the other for negative denominators. This, after all, is how practitioners deal with such ratios.

4. How to Account for Non-Proportionality in Ratios

Measurement using ratios requires proportionality between components. If the natural relationship between ratio components x and y is of the form $y = ax + b$ (non-proportional) rather than $y = ax$ (proportional), then the measurement will be misleading as x cannot have constant standards or norms (see, e.g., Lev and Suder, 1979). The existence of non-proportional ratios is established by Sudersanam and Taaffe (1995) and by others.

How to overcome problems posed by non-proportionality? The Law of Proportionate Effect acknowledges that changes dy/y (%) may be proportional, not to y itself, but to $y + \delta$. In this case, the natural form governing the generation of reported magnitudes will be:

$$\frac{\partial}{\partial \delta} \ln L = \frac{1}{\delta} \sum_{i=1}^n \frac{y_i}{y_i + \delta} \quad (17)$$

instead of (14) in turn δ may be either constant or size related. Constant δ stem $\delta = \rho$ from the effect of fixed costs in a cost structure. Indeed, in a time series, the existence of fixed costs leads to a constant displacement in the distribution of Operating Costs. Size related δ arise in similar cases (see discussion).

When δ is constant, conditions leading to lognormal (or in (10) generate in (24) distributions known as three-parametric lognormal with threshold δ (Archison and Breen, 1985). In econometric, lognormal distributions functions are simply the result of displacing lognormal distributions by δ (Figure 3).

Figure 3: The lognormal (solid line) and the three-parametric distributions (dashed line).



In order to generalize and non proportionality, one of the following ratios,

$$\frac{y_i}{y_i + \delta} \quad \text{or} \quad \frac{y_i - \delta}{y_i}$$

should be used for, respectively, three-parametric lognormal in the case of dependence ratios. Thresholds are δ or $\delta - y$. Indeed, there exists a ratio of y and δ for which the above non-proportionality ratios are the proportion 1.

How choosing thresholds are for the size measurement? Due to the exponential character of y , reported magnitudes may attain values many times larger than the threshold and, in such case, $y > \delta > y$. Non proportionality is significant only where magnitudes are not much larger than δ . For example, in the time series context, indeed, comparatively large size-independent thresholds are plausible only in a cross context. Observations in a time series have their origin in the same object, one firm or different

generally. Values that are constant inside firms, such as fixed costs, may create comparatively large thresholds. In cross-section, as observations have their origin in different objects, size-independent thresholds would require the existence of industry-wide "fixed" costs. Since such costs should allow for the survival of small firms, they must be small. An industry-wide cost of £4m for food manufacturers in the UK would be only 0.2% of United Biscuits' reserves, but it would equal or exceed the turnover of the 8% smallest firms in the industry.

As mentioned, fixed costs are also likely to generate size-related thresholds, mainly in cross-section where large firms have large fixed costs and small firms have small fixed costs. In this case, the behaviour of δ is similar to that of any accounting variable where correlation with size is the rule. According to (2), τ will be larger than expected for comparatively large δ (e.g., the case of large firms) and τ is not constant. As a consequence, size-related thresholds do distort ratio measurement. Notice that the problem, in this case, is not any displacement in the distribution of ratio components. Since δ are small for small v and large for large v , distributions are not displaced. The problem is non-linearity in their relationship. Indeed, size-related thresholds require the use of ratios of the type

$$\frac{v}{v + \delta}$$

For a specific value of β , the ratio will have a constant standard thus allowing measurement. On a logarithmic scale, the functional form of such measurement is

$$\log v = \beta \log v + \rho + \epsilon \quad (3)$$

where ρ is the logarithm of the ratio standard and ϵ is the observed deviation from that standard. (3) is similar to a regression. The slope β is approximate to the unit in the case of strict proportionality. Slopes smaller than 1 mean negative δ . In cross-section, they bias large firm's ratios downwards, mimicking scale effects.

How about the two types of δ just outlined (constant and size-related); be estimated? The ratio Fixed Assets (FA) to Current Assets (CA) is now used in a cross-section example. Five models of ratio, as follows, are compared. Each model is presented together with its logarithmic counterpart as multiplicative formulations require logarithmic scaling prior to coefficient estimation. Figure 4 shows, on a logarithmic scale, how the usual ratio (solid line) compares with each model. Figure 5 reproduces figure 4 on the original scale (only the region near the origin is displayed).

Figure 4

Figure 4: The usual ratio (solid line) compared, on a logarithmic scale, with the slope ratio (Model 2), threshold ratios (Models 3 and 4), and the threshold plus slope ratio (Model 5).

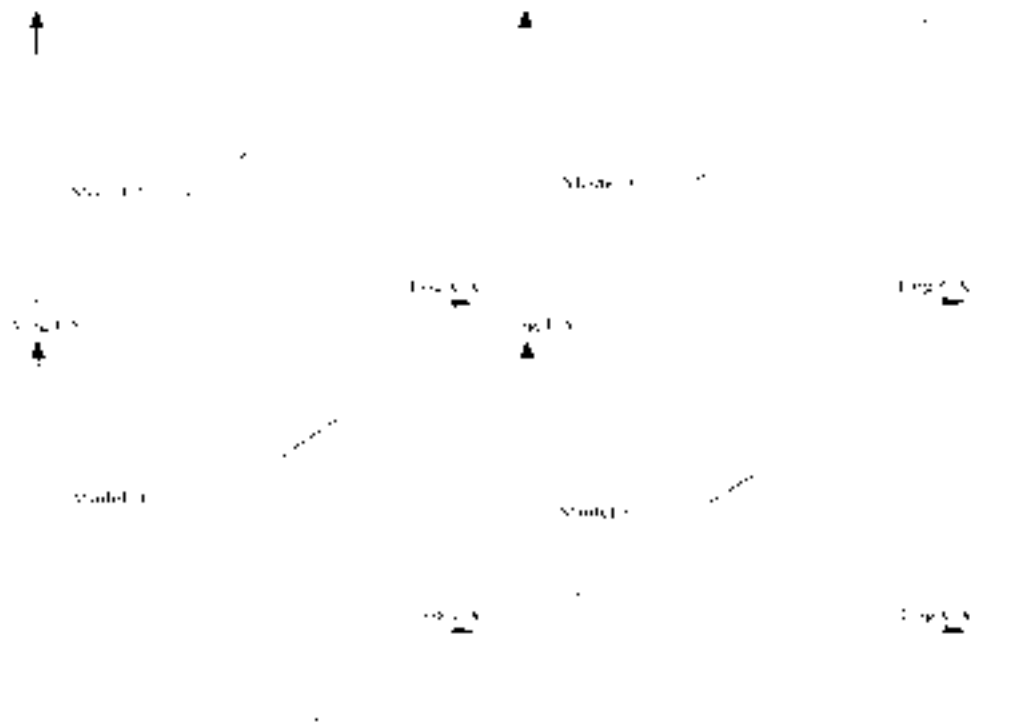


Figure 5: The usual ratio (solid line) compared with the slope ratio (Model 2), threshold ratios (Models 3 and 4), and the threshold plus slope ratio (Model 5).



Model 1: the usual ratio, no correction introduced. This ratio requires the estimation of one parameter, the standard. In cross-section, an appropriate standard is the median of the distribution of the ratio on a logarithmic scale, μ . Therefore

$$\frac{FA}{CA} = \exp(\omega + \log FA - \mu - \log CA) =$$

with ω being a random disturbance and $\omega \sim U$. When the distribution of CA and FA are nearly lognormal, the standard may be estimated by subtracting the averages of $\log FA$ and $\log CA$ (however, since a significant threshold in the distribution of CA is removed, such standard will be misleading).

Model 2: ratio with size-related threshold in the denominator. This ratio may control for the existence of economies of scale and other forms of non-linearity. It requires the estimation of two parameters, the standard μ (on a logarithmic scale) and the slope β . Regression's standard (1/5) may be used to estimate both.

$$\log(Y/A) = \mu + \delta \log(XA) + \epsilon \quad \text{corresponding to the ratio } \frac{Y/A}{XA^{\delta}} = 10^{\mu + \epsilon}$$

Graphically, the ratio is a straight line on logarithmic scale and a straight line on ordinary scale.

Model 3: ratio with constant threshold in the denominator (joint estimation). This ratio corrects for a constant threshold but in such a way that the standard is corrected as well. It requires the estimation of two parameters, the standard μ and logarithm δ , and the constant threshold δ . Both are estimated using

$$\log(YA) = \mu + \delta \log(XA + \delta) + \epsilon \quad \text{corresponding to the ratio } \frac{Y/A}{XA + \delta} = 10^{\mu + \epsilon}$$

Graphically, the ratio is a straight line on logarithmic scale and a straight line on ordinary scale.

Model 4: ratio with constant threshold in the denominator (independent estimation). In this case there is only one degree of freedom, μ , the logarithm of the ratio standard. The constant threshold δ of the distribution of $(Y/A) + \delta$ is estimated prior to that of the standard μ as in the next section.

$$\log(YA) = \mu + \delta \log(XA + \delta) + \epsilon \quad \text{corresponding to the ratio } \frac{Y/A}{XA + \delta} = 10^{\mu + \epsilon}$$

Graphically, the ratio is non-linear data and it will fit model 3 or logarithmic size and diameter on ordinary scale. This ratio probably is the most tested as non-proportionality is corrected for. It fits converging with the ratio on ordinary and normalized straight lines.

Model 5: ratio with both constant and size-independent thresholds. This ratio requires the estimation of three parameters, μ , δ and δ , using

$$\log(YA) = \mu + \delta \log(XA + \delta) + \epsilon \quad \text{corresponding to } \frac{Y/A}{XA + \delta} = 10^{\mu + \epsilon}$$

Graphically, it is similar to that of model 3 but in 5.

The independent estimation of a constant threshold (model 4) may be carried out using any of the procedures to detect the re-parametric lognormality in distributions. For example, $\delta = 0.00000$ is estimated using the procedure suggested by Royston (1982). The same estimation of δ and the ratio standard (model 3) and finally be carried out using more sophisticated algorithms.

Table 2 shows the variability explained (R²) by each of the five ratios. Skewness and kurtosis of residuals on a logarithmic scale (i.e., as displayed) As can be seen, by allowing β into model 2, R² approaches the variability explained by constant thresholds models (i.e., 1). Once δ is accounted for, β returns to its estimated value of nearly 1 (model 5).

Table 2: Parameters and statistics of the five models.

Model	α	β	δ	R ²	Skewness	Kurtosis
1	0.00	0.00	0.00	0.77	0.24	1.18
2	0.01	0.94	0.00	0.80	0.00	1.07
3	0.00	0.00	0.28, 0.01	0.81	0.11	1.16
4	0.16	0.00	0.28, 0.01	0.81	0.07	1.18
5	0.00	0.96	0.28, 0.01	0.81	0.12	1.18

Another example of the adequacy of the field ratios may be taken from practice: time-series regressions where Operating Costs explains Sales are traditionally used to estimate Fixed Costs and intercept term. Where the threshold ratios (and as an alternative to such regressions, the obtained) scenarios are clearly different in meaning. Where the correlation between Sales and Operating Costs is not very high, then the slope of the regression is clearly smaller than the ratio standard. In the limit, for a correlation approaching zero, such slope would also become zero and the intercept, which is supposed to estimate of Fixed Costs, would equal the expected value of operating costs. Regressions are thus inadequate for this task. They introduce in the estimation the spurious effect of correlation; by contrast, the threshold ratio correctly explains Fixed Costs as a displacement in the distribution of Operating costs.

5. How to Model Firm Size

Firm size is often chosen as a substitute for numerous theoretical constructs ranging from risk to liquidity or even political costs. Size is also an ingredient of its own primary theoretical models. In spite of this widespread use, size has remained a poorly defined concept. Where the use of size is required by theory, empirical studies typically revert to as the proxies such as Total Assets, Market Capitalization or Sales (see, e.g., Brock and Richardson, 1997).

The multiplicative character of magnitudes reported in accounts suggests, as stressed before, that the generation of their distribution is driven by size. It is indeed possible to derive a simple and effective definition of firm size from the two foregoing assumptions: first, reported magnitudes m indicate the effects Law of Proportionate Effect (see 01); financial ratios do remove the effect of size.

Implications

In order for ratios to be effective, the likelihood of observed discrepancies in relation to the standard must be independent of size. For instance, the Return on Assets ratio is useless if an increase or decrease of 2% has different meanings for small and large firms. The probability distribution of ratios must therefore be homoscedastic in terms of size. This implies, in the general case, that percent changes in both the numerator and denominator must be size independent. As mentioned, variables where percent changes are size independent are said to obey the Coburn's Law of Proportional Effect. Indeed, the widespread use of ratios agrees with the fact that reported numbers are multiplicative.

It was also mentioned that multiplicative mechanisms lead to broadly lognormal distributions. In fact, observed y generated as in (1) may be described functionally as

$$y = X(1 + \lambda)^t \quad \text{for discrete } t \text{ (or } y = X \exp \lambda t + \epsilon) \quad \text{for continuous } t \quad (4)$$

where λ is the variable which drives changes in y , X is a logarithmic expectation and ϵ is a random increment. The level X is the value of y for $t = 0$. It may be demonstrated that, any where cost nucleus computing is assumed, as in the right-hand side of (4), any ratios be validly used?

As for the second assumption mentioned above, that of ratio y/x eliminating the effect of size, it leads to the requirement that the numerator y and denominator x should both be generated under the effect of equal rates of change ($y = x$). In fact, where $y = x$ the ratio standard would not be constant, showing a rate of change of y/x with t . This requirement importantly suggests that, not just y and x but the other numbers found in a specific annual report are generated under the effect of the same rate of change. Indeed, the validity of the ratio method rests on the validity of several, widely used ratios (not just on one or two cases), where numbers used to form one of such ratios are also used to form other ratios. Therefore, if ratios are to be of any use, there must exist a common source of variability underlying numbers reported in the accounts of a firm in a given year. Where numbers x_1, x_2, \dots, x_n all belong to a specific annual report, this assumption leads to $x_1 = x_2 = \dots = x_n$. It is possible to show that such unique source of variability possesses the attributes of size.

How can size be estimated from this rate of change? A specific annual report, say, report i , is characterised by what value τ assumes x_l , the magnitude reported by item l , is explained as

$$x_l = X_l \exp[\lambda \tau + \epsilon]$$

where ν_j is the effect of size (the same for all items reported in j). In an additive form,

$$\log x_k = \mu_k + \sigma_j + \varepsilon \quad (5)$$

where $\mu_k = \log X_k$ and $\sigma_j = \nu_j$. Formulation (5) is basically an Analysis of Variance, i.e., a type of linear model aimed at explaining variability in terms of membership of discrete classes. Specifically, in (5) $\log x_k$ is explained by its membership of two classes, the item class, μ_k , and the annual report class, σ_j . The item class is a fixed (deterministic) effect, as it denotes the fact that k is a specific item amongst those in the sets of accounts reported by firms. These accounts are indeed fixed in number and in type. By contrast, the annual report's class is a random effect: it denotes the fact that j is one of the randomly selected annual reports in the sample. Each of these two classes possesses levels, namely, there can be as many levels of k as items in the sets of accounts; there can be as many levels of j as different reports in the sample.

In cross section, μ_k is the expected value of logarithmic magnitudes reported in item k , estimated as the mean of k calculated using all the annual reports in the sample. Thus, in this case, X_k in (5) is the median of magnitudes reported in k . The size effect, σ_j , is the expected $\log x_k - \mu_k$ for numbers in j and its estimation is straightforward: given N numbers, all of them reported in j , (5) is first applied to each of these numbers. The N formulations obtained are then added. Since σ_j is the same in all of these formulations, it is possible to write

$$\sigma_j = \frac{1}{N} \sum (\log x_k - \mu_k) + \frac{1}{N} \sum \varepsilon$$

Any source of variability common to all $\log x_k$ in (5), by construction, accounted for by σ_j . Therefore, even where correlation amongst some ε may exist, the term

$$\frac{1}{N} \sum \varepsilon$$

should tend to zero with an increasing N , leading to

$$\text{Estimate of } \sigma_j = \frac{1}{N} \sum (\log x_k - \mu_k)$$

An estimated size may thus be obtained simply by averaging the logarithms of appropriately adjusted magnitudes drawn from the firm's actual report. Exact confidence intervals for σ_j can also be obtained, the corresponding standard errors being t -distributed with $N-1$ degrees of freedom.

The estimation of σ_j faces two obvious difficulties. First, x_k cannot be drawn from all possible accounts because items such as Earnings, being a subtraction, may take on negative values and cannot be transformed into logarithms. This pre-selection of items introduces a bias in the estimation of size. Second, inside each annual report, some x_k are correlated. This increases the standard error. In practice, however, size can be estimated with accuracy by averaging the logarithms of positive-only magnitudes such as Cash and

Short-Term Investments, Receivables, Total Inventions, Property, Plant and Equipment, Debt, Common Equity, Number of Employees, Net Sales, Cost of Goods Sold, Research and Development, Depreciation, Depreciation of the year, Interest Expense, and others.

6. Concluding Remarks

Accounting data obey a set of simple rules. The first rule states that since reported numbers and ratios are multiplicative, logarithmic transformations should be used to bring data to additive behavior. The second rule explains that accounting identities often distort otherwise multiplicative distributions of ratios into unexpected shapes. Such identities should be accounted for prior to logarithmic transformation. The third rule explains the different ways to account for non-proportionality of errors. The fourth rule states that numbers in a specific report are generated under the influence of size, being possible to obtain a size estimate just by averaging several of these numbers conveniently adjusted. The distribution of numbers and the development of two-dimensional plots for analysis were also discussed.

It may be asked why so many contributions to the literature have led to a pessimistic view of ratios and accounting data in general. Reasons seem to lie in an apparent lack of theoretical drive, probably led by an *ad hoc* conviction that accounting data should be complex and full of exceptions, just as the production of such data is indeed complex and full of exceptions. It is expected that by using the results from this note, researchers and practitioners may find that such a *conviction* is rather a tradition with no scientific support.

References

Anderson, J. and Brown, J. (1977), *The Lognormal Distribution*, (Cambridge University Press).

Leenders, R. (1977), 'Pitfalls in the Application of Discriminant Analysis in Business Finance and Economics', *The Journal of Finance*, Vol. 32, No. 3, pp. 875–890.

Lev, B. and Srajer, S. (1998), 'Methodological Issues in the Use of Financial Ratios', *Journal of Accounting and Economics*, December, pp. 87–101.

Meloy, S. (1986), 'The Rate of Mean, the Mean of Ratios, and Other Benchmarks', *Journal of the Faculty of Finance Society*, Vol. 7, No. 1, pp. 73–90.

Meloy, S. and Ingstrup, D. (2002), 'Proportional Growth and the Theoretical Foundations of Financial Ratios', *Finance*, Vol. XXXIII, No. 3, pp. 397–416.

Royner, J. (1982), 'An Extension of the Shapiro and Wilk Test for Normality to Large Samples', *Applied Statistics*, Vol. 31, No. 7, pp. 115–122.

Sa, C. (1983), 'Some Empirical Evidence on Outliers and the Non-Normal Distribution of Financial Ratios', *Journal of Business Finance & Accounting*, Vol. 10, No. 1, pp. 183–196.

Srinivasan, P. and Taffler, R. (1995), 'Financial Ratio Proportionality and Intra-Temporal Stability: An Empirical Analysis', *Journal of Business and Finance*, Vol. XIX, pp. 48–69.

Tepfert, M. (1996), 'A Guided Theory of Financial Ratios', *Accounting and Business Research*, Vol. 26, No. 51, pp. 17–28.

Ingstrup, D. (1995), 'Accounting Identities and the Distribution of Ratios', *British Accounting Review*, Vol. 27, No. 2, pp. 99–112.

Notes

1. Natural firms do not behave mechanistically rather than observations. They are often expressed as differential or as difference equations.
2. Eisenberg (1977) mistakenly stated that 'logarithm-formed variables give less weight to equal percentage changes in a variable when the values are large than when they are small'. The implication would be that one does not believe that there is a much difference between a \$1 billion and a \$2 billion size firms as there is between a \$1 million and a \$2 million size firms. The percentage difference in the log will be greater in the latter than in the former case (up 50%). Eisenberg's point is that the logarithm of proportions of the log transformed measurement is equivalent to comparing proportions twice. This inappropriate warning against logarithmic transformation gave support of the true to the use of *z*-test techniques such as those proposed by Fama and French (1987).

- 3. The coefficient of variation is the standard deviation expressed as a fraction of expected value
- 4. Royston uses trial and error to find out which δ maximises the parameter W of the Shapiro and Wilk's test of normality
- 5. (1) may not necessarily lead to (4). Indeed, the simplest formulation generated by (1) is

$$x = V_0 \exp[\delta x - \sigma^2 \delta^2 x^2 / 2]$$

rather than (4). Since the compounding effect is now influenced by the standard deviation of z , not just by size, this variable is non-proportional and cannot be used to form ratios. In case, however, of continuously compounding rates of change, then a proportional mechanism is obtained.