# EFFECTS ON TREND STATISTICS OF THE USE OF MULTIPLICATIVE NOISE FOR DISCLOSURE LIMITATION

B. Timothy Evans, Bureau of the Census
Statistical Research Division, Rm. 3224-4, Bureau of the Census, Washington, DC 20233 [1]

Key Words: disclosure limitation, confidentiality, multiplicative noise, trend statistics, tabular data

## 1. Background

Historically the Census Bureau has favored disclosure limitation methods that protect sensitive data by limiting the amount of information given out. However, the Bureau is now considering methods that would allow for the release of more information but at the cost of having to distort the data in some way (Zayatz, Moore, and Evans, 1996). In the case of establishment tabular data, the traditional approach has been to suppress the publication of cells that are deemed sensitive, i.e., at risk for disclosing an individual respondent's data. Other cells, called complementary suppressions, must then also be suppressed to prevent the values of sensitive cells from being recovered through addition and subtraction of published cells. (For a complete discussion of cell suppression, see e.g., Federal Committee on Statistical Methodology, 1994.) Cell suppression thus protects sensitive data by limiting the amount of information given in the tables.

Cell suppression has its disadvantages, however. It withholds information that is not sensitive, namely the complementary suppressions. The process of choosing complementary suppressions is a complicated and time-consuming operation. And suppression patterns must be coordinated among all tables; that is, if a cell is suppressed in one table then it must be suppressed in all other tables in which it appears. This last requirement creates tremendous difficulty in the fulfillment of requests for special tabulations following publication of standard tables.

In an effort to simplify the disclosure review process and to increase the amount of data that can be released, the Census Bureau has recently begun looking at alternatives to cell suppression for performing disclosure limitation on establishment tabular data. Thus far the research has focused on introducing noise into the establishment microdata records prior to tabulation. Noise addition would allow more cells to be published because it eliminates the need for complementary suppressions; sensitive cells are protected simply by the noise present in their published values. Also, noise would greatly simplify the disclosure limitation process because the noise only needs to be added once, and then any number of tabulations can be produced from the perturbed microdata. There would be no worries about consistency of cell values between tables or about coordinating suppression patterns among all data products.

While using noise would allow for the release of more data, questions remain about the usefulness of data that has been perturbed. Others have explored this question regarding the possibility of releasing perturbed economic microdata files (e.g., McGuckin and Nguyen, 1990), but to date little work has been done on the usefulness of tabular data in the presence of noise. Evans, Zayatz, and Slanta (1996) experimented with introducing multiplicative noise into establishment microdata prior to tabulation and found that resulting level estimates were generally not adversely affected. Muralidhar, Batra, and Kirs (1995) looked at descriptive statistics of distributions in statistical databases and found that adding noise to microdata provided sufficient security while preserving the accuracy of the descriptive statistics. They also observed that using multiplicative noise produced more useful data than using additive noise.

Observing the behavior of simple level estimates in the presence of noise is only the necessary first step, however. Data users use these level estimates to perform many types of analyses, such as describing relationships among data items or looking at the behavior of certain variables over time. It remains to be seen what effect noise may have on these analyses. This paper begins to address this issue by investigating the effects of noise on a simple type of analysis: year-to-year trends.

## 2. Formulation of the Problem

In assessing the effect of noise on a trend statistic, we will look at the ratio of the noisy trend to the true trend. Let $Y_1$ and $Y_2$ be the true (noise-free) level estimates of some variable Y for year 1 and year 2 (not necessarily consecutive), respectively. Let $R = \frac{Y_2}{Y_1}$. The true trend in Y (expressed as a decimal rather than a percent) between the 2 years is $\frac{Y_2 - Y_1}{Y_1}$, which is equal to

---

R - 1.

When noise is added to the underlying microdata, all estimates produced from that microdata will contain at least a small amount of noise. The noise derives from individual multipliers being applied to individual observations and then being summed. For simplicity, assume we are dealing with unweighted data, and express $Y_1$ as $Y_1 = \sum_i y_{1,i}$, where $y_{1,i}$ is establishment i's value of Y in year 1. Then the noise-added estimate of Y in year 1 can be written as

$$\text{noisy } Y_1 = M_1 Y_1 = \sum_i m_{1,i} * y_{1,i} \, ,$$

where $m_{1,i}$ is the multiplier associated with establishment i in year 1 and $M_1$ can be described as the *net* noise multiplier for $Y_1$. Explicitly, $M_1 = \frac{\text{noisy } Y_1}{\text{true } Y_1}$. Similarly, let $M_2$ be the net multiplier for year 2. Note that $M_1$ and $M_2$ are not known in advance. Assuming the strategy described in Evans, Zayatz, and Slanta (1996) for assigning the values of the $m_i$'s, $M_1$ and $M_2$ will generally be much closer to 1 than the individual $m_i$'s, and their distance from 1 will depend on the skewness of the distribution of the $y_i$'s. (The assignment scheme recognizes and attempts to accommodate the fact that most economic data distributions are inherently skewed.) The more skewed the distribution of the establishments contributing to a particular cell estimate, the less likely it is that noise in individual establishments will cancel out as establishments are aggregated, and hence the farther from 1 we would expect the net noise multipliers to be.

Note also that, since the establishment multipliers were selected such that $E(m_i) = 1$, the expected value (given the $y_i$'s) of the resulting net noise multiplier is also 1:

$$E\,(M_1) = E\left(\frac{\text{noisy } Y_1}{\text{true } Y_1}\right) = \frac{1}{\text{true } Y_1} * E\left(\sum_i m_{1,i} * y_{1,i}\right)$$

$$= \frac{1}{\text{true } Y_1} * \sum_i E(\,m_{1,i}\,) * y_{1,i}$$

$$= \frac{1}{\text{true } Y_1} * \sum_i y_{1,i} = 1$$

We are interested in how the noise in the component level estimates translates into noise in the trend. Our measure of the noise in the trend is the trend's net noise multiplier, which we will denote $M_{\text{trend}}$. Specifically, $M_{\text{trend}} = \frac{\text{noisy trend}}{\text{true trend}}$. The noisy trend is the trend computed using the noise-added level estimates and can be written as $\frac{M_2 Y_2 - M_1 Y_1}{M_1 Y_1}$. We can then express $M_{\text{trend}}$ as follows:

$$M_{\text{trend}} = \frac{\frac{M_2 Y_2 - M_1 Y_1}{M_1 Y_1}}{\frac{Y_2 - Y_1}{Y_1}} = \frac{\left(\frac{M_2}{M_1}\right)\frac{Y_2}{Y_1} - 1}{\frac{Y_2 - Y_1}{Y_1}} = \frac{\left(\frac{M_2}{M_1}\right)R - 1}{R - 1}$$

The amount of noise resulting in the trend thus depends on two quantities. First, it depends on R, the true ratio of the Y values between the 2 years. It stands to reason that if there is little change in a variable Y between 2 years, then the year-to-year change would be very easily obscured by even a small amount of added noise. If the change in Y is very small, that is if $Y_2$ is very nearly equal to $Y_1$, then R is close to 1 and $M_{\text{trend}}$ will have a tendency to be very large (in magnitude). Note in particular that if $R = 1$, then $M_{\text{trend}}$ is undefined; in this case it is impossible to express any noise in the trend as a percentage of the true value because the true trend is 0.

Secondly, the amount of noise in the trend depends on $\frac{M_2}{M_1}$. Regardless of the magnitude of the true trend, notice that the closer $M_2$ is to $M_1$, the closer $M_{\text{trend}}$ will be to 1. That is, if the values of Y in the 2 years ended up with the same amount of noise in them, the common net noise factor would cancel when computing the trend $\left(\frac{M_2\,Y_2 - M_1\,Y_1}{M_1\,Y_1} = \frac{M_{\text{common}}\,(\,Y_2 - Y_1)}{M_{\text{common}}\,Y_1} = \frac{Y_2 - Y_1}{Y_1}\right)$ and the resulting trend would have no noise in it at all.

It is worth noting that $M_{\text{trend}}$ does not depend directly on the values of $M_1$ and $M_2$ individually, only on their relative sizes. Even if $M_1$ and $M_2$ are both very far from 1 (as in a single-contributor cell), if they are far from 1 in the same direction and by about the same amount, the trend can still end up with almost no noise in it. Two very noisy level estimates do not necessarily produce a noisy trend.

Values of R close to 1 will tend to make $M_{\text{trend}}$ large in magnitude (in either the positive or negative direction, depending on whether $M_2$ is larger or smaller than $M_1$), while values of $\frac{M_2}{M_1}$ close to 1 will tend to bring $M_{\text{trend}}$ closer to 1. Which is the stronger force?

## 3. Amount of Noise in Trends

The first question we would like to answer regarding trends is whether the addition of noise to the microdata results in any bias in the trends computed from the noisy level estimates. We have already seen in Section 2 that the level estimates themselves are unbiased $[E(M_1) = 1]$; do unbiased level estimates result in an unbiased trend? Using the Taylor series expansion result that $E(\frac{X}{Y}) \approx \frac{E(X)}{E(Y)}$ and treating R as fixed, we see that

$$E(M_{\text{trend}}) = E\left(\frac{\frac{M_2}{M_1}R - 1}{R - 1}\right) = \frac{R}{R - 1} * E\left(\frac{M_2}{M_1}\right) - \frac{1}{R - 1}$$

$$\approx \frac{R}{R - 1} * \frac{E(M_2)}{E(M_1)} - \frac{1}{R - 1} = \frac{R}{R - 1} - \frac{1}{R - 1} = 1$$

In order to verify this unbiasedness, and to assess the amount of noise that will typically be present in a trend (since in individual applications $M_{\text{trend}}$ will in

general not be 1), we conducted an experiment using 4 years' worth (1990-1993) of data from the Census Bureau's County Business Patterns (CBP). For three variables and for a number of 2-digit SIC (Standard Industrial Classification) codes, we added noise to the microdata and computed trends for each cell and for each pair of years. This resulted in a total of 3486 trends. We replicated the addition of noise and computation of trends 100 times and observed the behavior of the trends over all replications.

For each trend, we computed the average value of $M_{trend}$ over all replications and looked at the distribution of this quantity over all trends. The distribution was centered exactly at 1, with very narrow spread, thus bearing out the theoretical result.

The next question we would like to answer is how much noise we can typically expect to be present in a trend. The amount of noise in the trend can be measured several ways. One way is to look at the noise relative to the size of the original trend, i.e., as a percent of the percent change. Using the CBP data, for each cell and for each replication we computed the relative percent noise in the trend as $\left| \frac{\text{noisy trend - true trend}}{\text{true trend}} \right| * 100\%$. (Note that this is equal to the absolute value of $M_{trend} - 1$, expressed as a percent.) For each cell, we looked at the distribution of this quantity over all 100 replications and computed selected percentiles. Then we looked at the distributions of these percentiles over all trends. So we are looking at distributions, over all *trends,* of percentiles which are themselves computed from distributions, over all *replications*, of the absolute relative noise present in an individual trend.

Among all trends with a true value larger than 1 percent in magnitude (i.e., true trend < -1% or > +1%), the distribution of Q1 as described above had quartiles Q1 = 2.06% and Q3 = 11.31%. The distribution of Q3 had quartiles Q1 = 5.00% and Q3 = 19.85%. Even the maximum Q3 over all trends was 113.78%, meaning that even in the worst cases it only happened slightly more than 25% of the time that the noise-added trend was more than twice as large (in magnitude) as the noise-free trend. Thus, judging by the behavior of the quartiles over all cells, we can typically expect noise-added trends to contain about 2 to 20 percent noise, relative to the true value of the trend.

It was generally true, moreover, that the larger values of the replication-distribution quartiles corresponded to trends whose true values were small. In fact, it was this tendency for small trends to contain large amounts of *relative* noise that led us to exclude very small trends (smaller than 1 percent in magnitude) from the preceding analysis. This phenomenon is not surprising and illustrates the limitations of looking at percent changes in a statistic that is itself a percent change rather than an actual quantity. Because most percent change statistics tend to be relatively small, changes in the magnitude of the percent change that are small in absolute terms appear very large when viewed as a fraction of the statistic itself.

As a more familiar example, consider that even a moderate-sized year-to-year change of 4.4%, for instance, has a substantial amount of "noise" introduced into it (relative to the size of the change) simply by being rounded to the nearest whole percent. Expressing this trend as .04 rather than .044 has changed its value by 9%, a percentage that would be highly objectionable in a level estimate.

Recognizing the limitations of expressing the noise in trends in relative terms, we also looked at the noise in absolute terms. For each trend and for each of the 100 replications, we computed the absolute difference in percentage points between the noisy trend and the true trend, i.e., $\left| \text{noisy trend} - \text{true trend} \right| * 100\%$. (For example, if the true trend was .02 and the noise-added trend was .04, we would describe this as a 2% difference in this instance, as opposed to a 100% difference in relative terms.) We then averaged this quantity over all replications for each trend and looked at the distribution of the average over all trends.

Using the average absolute percentage point difference as the measure, the median amount of noise over all trends was only 0.7%, and even the 90th percentile was only 2.89%. In particular, among trends having a true value of less than 1% in magnitude, the median amount of noise was only 0.5%. In relative terms this amount would appear very large but is in fact only on the order of rounding error when looked at in absolute percentage points. This time, points in the right tail of the distribution tended (unsurprisingly) to correspond to very large true trends, for which these large absolute differences would appear small in relative terms.

To summarize, the extent to which adding noise to the underlying microdata (and thence to the level estimates used in computing the trend) introduces noise into trend statistics depends on how the amount of noise in the trends is measured. When viewed in relative terms, the amount of noise in a trend will typically be in the range of 2 to 20 percent. The amount can potentially be much higher, but these higher values tend to correspond to small values of the true trend; in these cases measuring noise in relative terms makes the situation look worse than it is. When viewed as a straightforward difference between the noisy trend and the true trend, the amount of noise will typically be only 1 or 2 percent. In either case, whether these levels of noise in trends are acceptable is a question for further discussion and is beyond the scope of this paper.

## 4. Apparent Changes in Direction of Trend

However the amount of noise in a trend is measured, certainly a critical issue in reporting percent changes is whether the change is significantly different from 0. In this light, we would like to know under what conditions the presence of noise in level estimates might cause a trend to appear to change sign.

If $M_{trend} < 0$, this means that the addition of noise to the level estimates caused the trend to change sign. Intuitively we would expect that the true trend would have to be very small in order to be so adversely affected by the noise as to appear to change direction. Is this the case?

Note that $M_{trend} < 0 \Leftrightarrow \dfrac{\left(\frac{M_2}{M_1}\right)R - 1}{R - 1} < 0$. If $R > 1$, then

$M_{trend} < 0 \Leftrightarrow \left(\frac{M_2}{M_1}\right)R - 1 < 0 \Leftrightarrow \left(\frac{M_2}{M_1}\right)R < 1 \Leftrightarrow \frac{M_2}{M_1} < \frac{1}{R} < 1$. It stands to reason that this requires $M_2 < M_1$, considering that the true trend is upward but the noise makes it appear to be downward. If $R < 1$, then $M_{trend} < 0 \Leftrightarrow \frac{M_2}{M_1} > \frac{1}{R} > 1$ by a similar argument.

In either case, in order for $M_{trend}$ to be $< 0$, $\frac{M_2}{M_1}$ and its multiplicative inverse, taken as a pair, must be "farther away" from 1 than $R$ and $\frac{1}{R}$ are, but in opposite directions. More precisely, the interval $\left(\frac{1}{R}, R\right)$ [or $\left(R, \frac{1}{R}\right)$, depending on the size of $R$ relative to 1] must be contained in the interval $\left(\frac{M_2}{M_1}, \frac{M_1}{M_2}\right)$ [or $\left(\frac{M_1}{M_2}, \frac{M_2}{M_1}\right)$ if $R < 1$].

This condition is very rarely met, mainly because of the restrictions imposed in Evans, Zayatz, and Slanta on the updating of individual establishment noise multipliers from one period to another. These restrictions were designed to maintain the utility of trend statistics, and the net result is that $M_2$ seldom differs from $M_1$ by more than about 1%, whereas most trends are larger than this.

The results from the test with County Business Patterns data reinforce this assertion. Of the 3486 trends examined, only 376 (slightly more than 10%) changed sign *even once* over all 100 replications. Of these, about 50% were trends whose true values were less than 1% in magnitude, again illustrating the susceptibility of very small trends to being obscured by even a small amount of noise. Of the remaining 50%, the majority were trends in cells that were dominated by a very large contributor or contributors, i.e., sensitive cells. In such cells, the net noise multipliers $M_1$ and $M_2$ are very close to (and in single-contributor cells, identical to) the establishment-level multipliers $m_1$ and $m_2$ assigned to the dominant contributor. (If the cell is dominated by 2, 3, etc. contributors that are all perturbed in the same direction, the effect will be similar to that of a single dominant contributor.) It is not uncommon for these individual establishment multipliers to differ from each other by several percent, so for these sensitive cells that exhibit micro-level behavior, the ratio of the net noise multipliers can easily exceed the size of the true trend, resulting in the trend appearing to change direction. (Section 5 discusses sensitive cells in more detail.)

For the majority of cells, then, noise does not cause the trends to change direction. The only exceptions are trends that are very close to zero to begin with (in which case measurement errors probably make their true direction questionable anyway) and trends in sensitive cells, for which an apparent change in direction is not necessarily undesirable and can even be looked at as a form of protection.

## 5. Differing Effects on Sensitive vs. Nonsensitive Cells

A final area of concern is whether there will be differences in the amount of noise that typically results in trend statistics for sensitive cells as compared to nonsensitive cells. In assigning noise multipliers to establishments, the goal was to ensure that sensitive cells, whose values (i.e., level estimates) need to be protected, would receive large amounts of noise, while at the same time trying to minimize the amount of noise that would appear in cells that aren't at risk for disclosure. Ideally we would like the same to be true for trend statistics – that trends for nonsensitive cells remain relatively untouched by the addition of noise but that trends for sensitive cells be more noticeably distorted.

Examination of the County Business Patterns data yields mixed results in this regard. When noise is measured in relative terms, a comparison of the distributions of the amount of noise in sensitive cells vs. nonsensitive cells indicates that trends in sensitive cells get only slightly more noise than those in nonsensitive cells. When measured as a simple difference between noisy and noise-free trends, the amount of noise in sensitive trends is much greater than in nonsensitive trends, but even here there is some cause for reservation.

We observed that the most extreme (largest in magnitude) year-to-year changes almost always occur in sensitive cells. As mentioned in the previous section, values in sensitive cells reflect the behavior of only one or two dominant companies, while nonsensitive cells generally have many contributors and hence describe more aggregated, macro-level behavior. Naturally, we expect more variability at the micro level. Just as noise in individual establishments has a tendency to cancel out as the establishments are aggregated into cell totals, so too will highly divergent percent changes in individual establishments tend to produce a more moderate estimate of change as more establishments are added together.

Conversely, it is the sensitive cells whose values of R have the potential to differ the greatest from 1 by virtue of their describing what is effectively micro-level behavior. Considering the results of Section 3, these sensitive cells would be the most immune of all cells to having their trends disturbed by the addition of noise.

## 6.　　Conclusions

The situation regarding sensitive vs. nonsensitive trends raises a very important question: How much protection do we need to give to trend statistics? If we are looking at a cell that is dominated by one large contributor, is it sufficient to protect the level estimates and simply let the trends fall where they will, even if the trend in the presence of noise is virtually identical to the true trend? Would the noise in the level estimates discourage users from putting too much faith in the trend for that cell, even if in actuality the trend computed from the noisy estimates were a very close approximation to the true value?

Also at issue is how protection is measured, and more generally how the amount of noise is measured, for trend statistics. Measuring noise in relative terms can overstate the seriousness of the distortion in trends that were very close to zero before noise was added. On the other hand, measuring noise simply as the difference between the noisy trend and the true trend can imply more protection than actually exists for trends that were very large to begin with. It is our feeling that neither method is satisfactory and that better criteria need to be developed for evaluating the amount of noise present in trends.

Finally, as mentioned in Section 3, we need to decide how much noise is acceptable in trends for cells that aren't sensitive. Noise is desirable in sensitive trends as a means of protecting respondent data by preventing data users from recovering the true value of the trend from published results. At what point does the amount of noise in *nonsensitive* trends preclude users being able to draw meaningful conclusions from them as well?

The answers to the above questions will probably vary from one survey to another and from one data user to another. In any event, they cannot be answered mathematically (other than the measurement question) and must therefore be left to policy makers.

## 7.　　References

Evans, B.T., Zayatz, L., and Slanta, J. (1996), "Using Noise For Disclosure Limitation of Establishment Tabular Data," Proceedings of the 1996 Annual Research Conference, Washington, DC: U.S. Bureau of the Census.

Federal Committee on Statistical Methodology (1994), Statistical Policy Working Paper 22: Report on Statistical Disclosure Limitation Methodology, Washington, DC: U.S. Office of Management and Budget.

McGuckin, R.H. and Nguyen, S.V. (1990), "Public Use Microdata: Disclosure and Usefulness," Journal of Economic and Social Measurement, Vol. 16, pp. 19-39.

Muralidhar, K., Batra, D., and Kirs, P.J. (1995), "Accessibility, Security, and Accuracy in Statistical Databases: The Case for the Multiplicative Fixed Data Perturbation Approach," Management Science, Vol. 41, No.9, pp. 1549-1564.

Zayatz, L., Moore, R.A., and Evans, B.T. (1996), "New Directions in Disclosure Limitation at the Census Bureau," Proceedings of the Section on Government Statistics, American Statistical Association, pp. 89-93.