
SURVEYS AND CROSSOVERS

This section provides surveys of the literature in investment management or short papers exemplifying advances in finance that arise from the confluence with other fields. This section acknowledges current trends in technology, and the cross-disciplinary nature of the investment management business, while directing the reader to interesting and important recent work.

UNREALISTIC EXPECTATIONS: THE FUTILITY OF PRECISELY ESTIMATING A STOCK'S EXPECTED RETURN

Sanjiv R. Das^a and Daniel Ostrov^b

*...dedicated to the memory of Mark S. Joshi,
who worked to make results like these better known*

We reprise the result that even under the best circumstances, it is impossible to use observed return data for a stock to determine its expected return with any useful precision in a reasonable time frame. This is because the sample mean of returns, which is the best estimator for the expected return, has a large standard error. More specifically, the formula for this standard error is $\frac{\sigma}{\sqrt{T}}$, where σ is the stock's annual volatility and T is the number of years over which the returns are sampled. We note in particular that this standard error formula is not reduced by increasing the frequency of sampling the returns within the given time frame T .

As an example of the ramifications of this standard error formula, a stock with even a small volatility like 10% would still require 400 years of observed returns to determine (with 95% confidence) the stock's annual expected return within a percentage point, assuming the expected return is constant over these 400 years. Another important implication of this formula is that actual returns cannot be used to disprove any remotely reasonable estimate for the expected return.

In contrast, if the volatility of a stock or the covariance matrix of a number of stocks remains constant over time, they can be estimated from return information far more precisely within a reasonable time frame. Further, this precision can be improved by increasing the frequency of sampling the returns.

^a William and Janice Terry Professor of Finance and Data Science in the Leavey School of Business at Santa Clara University.

^b Professor of Mathematics and Computer Science at Santa Clara University.

The results of this paper are not new. For example, Appendix A of Merton (1980) also states the $\frac{\sigma}{\sqrt{T}}$ formula above. However, these results are also not that widely known or acknowledged. Because they are so fundamental for investment management and so easy to derive, we hope this short article with its simple exposition will reprise and remind the financial modeling community of the difficulty of working with time-series return data to estimate future returns.



1 Introduction

The three most common parameters of interest in modeling the evolution of stocks are the expected return of a stock, the volatility of a stock, and the covariance of any two stocks' returns. We can repeatedly sample returns for the stock (or stocks) and estimate these parameters. Specifically, we use the annualized sample mean of returns as the point estimator of the expected return, the annualized sample standard deviation as the point estimator of the volatility, and the annualized sample covariance as the point estimator of the covariance of two stocks' returns. The precision of these point estimators depends on their standard errors, meaning the standard deviations of the point estimators.

This paper derives these three standard errors in the best case scenario, where the model for stock evolution is as simple as possible and assumed to be completely accurate.

2 The Error in Approximating the Expected Return

This section describes the notation and the best case return model, then derives the standard error for the annualized sample mean using this model and its implications on computing a confidence interval for the annual expected return.

2.1 Notation and model for a single stock

Consider a single stock whose returns are sampled repeatedly, and the following notation:

- (1) T is the number of years over which returns are sampled.
- (2) h is the time step (in years) between samples.
- (3) $N = \frac{T}{h}$ is the number of samples.
- (4) $X_{h,i}$ is the i th sample return when the time step is h , where $i = 1, 2, \dots, N$.
- (5) Z_i are independent standard normal random variables, where $i = 1, 2, \dots, N$.
- (6) μ is the annualized expected return of the stock.
- (7) σ is the annualized volatility of the stock.
- (8) $\bar{X}_h = \frac{\sum_{i=1}^N X_{h,i}}{n}$, the sample mean of the $X_{h,i}$.
- (9) $\bar{X}_a = \frac{\bar{X}_h}{h}$ is the annualized sample mean.

With this notation, the goal is to determine how close the point estimator \bar{X}_a is to the parameter μ . This depends on the standard error, which is $SD[\bar{X}_a]$, the standard deviation of the point estimator \bar{X}_a .

Assume the *most simple* model for stock returns holds:

$$X_{h,i} = \mu h + \sigma \sqrt{h} Z_i, \quad (1)$$

where $i = 1, 2, \dots, N$. Note that this model includes geometric Brownian motion. Further,

assume that μ and σ are constant over the time frame T , no matter how long T is. This is the *best* possible case.

2.2 Calculating the standard error

Given this ideal model for stock returns, the point estimator \bar{X}_a for the parameter μ is

$$\begin{aligned}\bar{X}_a &= \frac{\bar{X}_h}{h} = \frac{\sum_{i=1}^N X_{h,i}}{hN} \\ &= \frac{\sum_{i=1}^N (\mu h + \sigma \sqrt{h} Z_i)}{hN} \\ &= \mu + \frac{\sigma}{\sqrt{hN}} \sum_{i=1}^N Z_i.\end{aligned}\quad (2)$$

Because the sum of normal random variables is normal, \bar{X}_a has a normal distribution. Taking the variance of the expression above and using that (i) the Z_i are independent of each other, (ii) $V[Z_i] = 1$, and (iii) $T = hN$ yields

$$\begin{aligned}V[\bar{X}_a] &= V\left[\mu + \frac{\sigma}{\sqrt{hN}} \sum_{i=1}^N Z_i\right] \\ &= \frac{\sigma^2}{hN^2} \sum_{i=1}^N V[Z_i] \\ &= \frac{\sigma^2}{hN} = \frac{\sigma^2}{T}.\end{aligned}$$

Finally, taking the square root gives the desired expression for the standard error:

$$\boxed{SD[\bar{X}_a] = \frac{\sigma}{\sqrt{T}}.}$$

Note that this expression does not depend on h . That is, increasing the frequency of sampling

within the time frame of T years does *not* improve the estimate of μ .

2.3 Confidence intervals for μ

Since \bar{X}_a has a normal distribution, the 95% confidence interval for μ is $\bar{X}_a \pm 2 \times SD[\bar{X}_a]$ —that is, $\bar{X}_a \pm \frac{2\sigma}{\sqrt{T}}$.

Assume for the sake of specificity in the following examples that $\bar{X}_a = 10\%$ over T years of returns (although the value of \bar{X}_a has no effect on the width of the confidence interval, of course). If $\sigma = 20\%$ and $T = 10$ years, the 95% confidence interval for μ is between -2.65% and 22.65% . This, of course, is an absurdly large confidence interval and of little use to financial planners. Even if the stock had an unusually small σ value of, say, 10% , the 95% confidence interval for μ would be from 3.68% to 16.32% . While this is half the size of the interval if $\sigma = 20\%$, it is still too large to be of real use to financial planners.

The only remaining lever to reduce the size of the confidence interval is to increase T . However, the confidence interval width depends on the *square root* of T , so to halve the confidence interval again so that it is from 6.83% to 13.16% , requires changing T from 10 years to 40 years. Indeed, to reduce the 95% confidence interval to be $\pm 1\%$ (so that for $\bar{X}_a = 10\%$, the confidence interval would be between 9% and 11%) would require sampling returns over $T = 400$ years. Of course, the assumption that the model for returns holds over such long time horizons, including μ and σ remaining constant, becomes questionable.

3 The Error in Approximating the Volatility

3.1 Notation and model

We retain the same model as before from equation (1) with μ and σ remaining constant. To the

previous notation, add the following:

- (1) S_h is the sample standard deviation of the $X_{h,i}$, where $i = 1, 2, \dots, N$.
- (2) $S_a = \frac{S_h}{\sqrt{h}}$ is the annualized sample standard deviation.

Given this notation, the goal is to determine how close the point estimator S_a is to the parameter σ . This means determining the standard error, which is $SD[S_a]$, the standard deviation of the point estimator S_a .

3.2 Calculating the standard error

The sample standard deviation S_a can be expressed by applying the model in equation (1) to $X_{h,i}$ and applying equation (2) to \bar{X}_h :

$$\begin{aligned}
 S_a &= \frac{S_h}{\sqrt{h}} = \frac{1}{\sqrt{h}} \sqrt{\frac{\sum_{i=1}^N (X_{h,i} - \bar{X}_h)^2}{N - 1}} \\
 &= \sigma \sqrt{\frac{\sum_{i=1}^N \left(Z_i - \frac{\sum_{i=1}^N Z_i}{N} \right)^2}{N - 1}} \\
 &= \sigma \sqrt{\frac{\left(\sum_{i=1}^N Z_i^2 \right) - \frac{\left(\sum_{i=1}^N Z_i \right)^2}{N}}{N - 1}}.
 \end{aligned}$$

Since the Z_i are independent, $\sum_{i=1}^N Z_i$ is a normal random variable with a variance of N , so $\frac{\left(\sum_{i=1}^N Z_i \right)^2}{N}$ has a Z^2 distribution. As N becomes large, this Z^2 term in the numerator becomes small compared to the $\sum_{i=1}^N Z_i^2$ term, and the 1 in the denominator becomes small compared to N . That is, as $N = \frac{T}{h}$ becomes large,

$$S_a \sim \sigma \sqrt{\frac{\sum_{i=1}^N Z_i^2}{N}}$$

$$\begin{aligned}
 &= \frac{\sigma}{\sqrt{N}} \sqrt{\sum_{i=1}^N Z_i^2} \\
 &= \sigma \sqrt{\frac{h}{T}} \sqrt{\chi_N^2} = \sigma \sqrt{\frac{h}{T}} \chi_N,
 \end{aligned}$$

where χ_k is the chi random variable with k degrees of freedom, which is defined as the square root of χ_k^2 , the chi-squared random variable with k degrees of freedom. The specific formula for the variance of the chi random variable¹ is $k - 2 \left(\frac{\Gamma(\frac{k+1}{2})}{\Gamma(\frac{k}{2})} \right)$ where $\Gamma(\cdot)$ is the gamma function, but the key feature of this variance formula is that it quickly converges to $\frac{1}{2}$ as k gets large, as seen in Figure 1. Therefore,

$$V[S_a] \sim \sigma^2 \frac{h}{T} V[\chi_N] \sim \sigma^2 \frac{h}{2T}.$$

Taking the square root of this expression gives the desired formula for the standard error:

$$SD[S_a] \sim \sigma \sqrt{\frac{h}{2T}},$$

as N , the number of samples, gets large.

In practice N does not have to be very large for this approximation to hold well. Being above 100 is certainly sufficient. Note that this expression for $SD[S_a]$, in contrast to the expression for $SD[\bar{X}_a]$ from the previous section, does depend on h , so it gets smaller as we increase the frequency of sampling within the time frame of T years.

3.3 Confidence intervals for σ

Since the chi distribution with 100 or more degrees of freedom is, essentially, normal (see Figure 2 for the case with 100 degrees of freedom), the 95% confidence interval for σ is, essentially, $S_a \pm 2 \times SD[S_a]$ —that is, $S_a \pm \sigma \sqrt{\frac{2h}{T}}$.

Consider the case where the sample standard deviation is $S_a = 20\%$, which is obtained by sampling

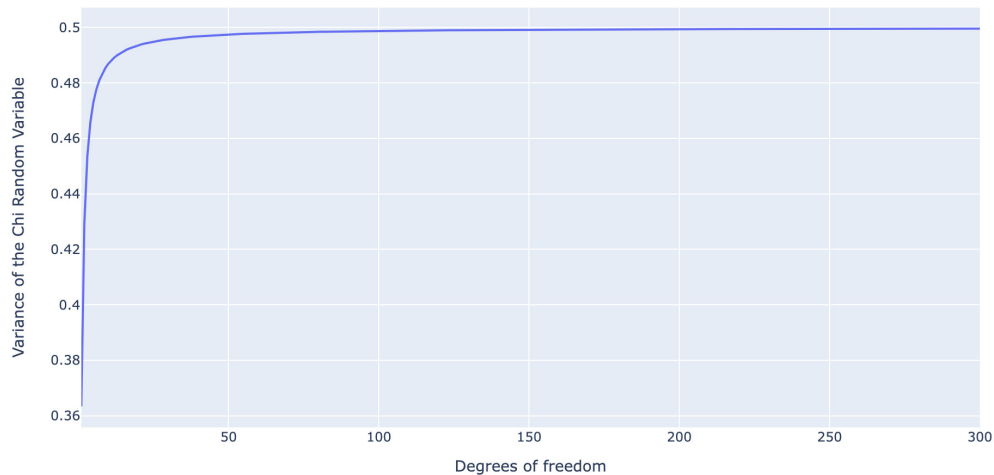


Figure 1 The variance of the chi random variable χ_k . As k , the degrees of freedom, increases, the variance quickly gets close to $\frac{1}{2}$. Indeed, if $k > 125$, the variance is between 0.499 and 0.5.

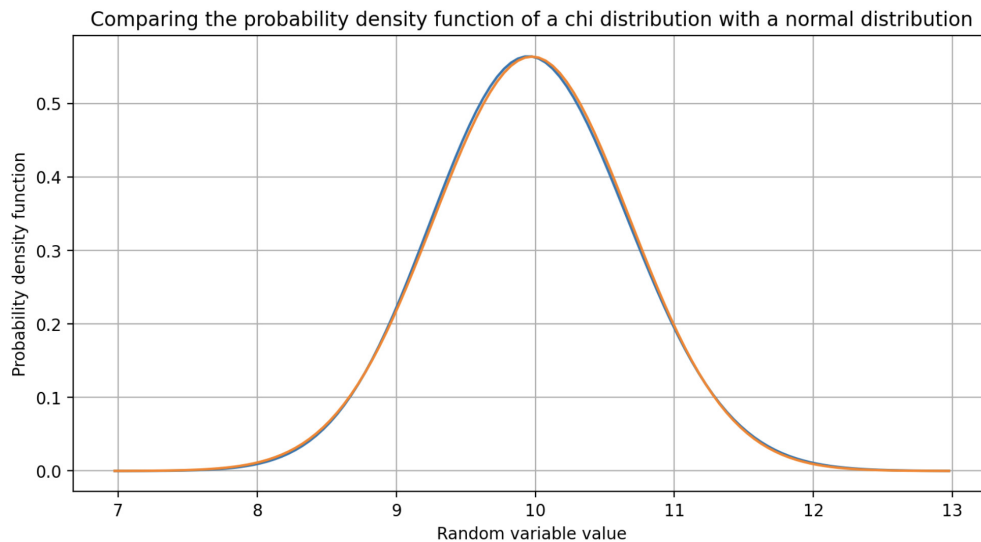


Figure 2 The distribution (in blue) for χ_{100} , the chi random variable with 100 degrees of freedom, compared to a normal distribution (in orange) with the same mean and standard deviation. They are nearly identical. This is also the case for chi random variables with more than 100 degrees of freedom.

the stock's returns four times in each of the 252 trading days over a year, so $h = \frac{1}{4 \times 252}$ and $T = 1$. Substituting S_a for σ in the confidence interval formula gives 95% confidence that σ is between 19.11% and 20.89%, with the small width of this interval justifying substituting S_a for σ . This delivers a very good estimate of σ for most purposes, and it can be made better, if needed, by sampling the returns more often each day.

4 The Error in Approximating the Covariance of Two Stocks' Returns

4.1 Notation and model

For multiple stocks, in addition to approximating the expected return and volatility of each stock, the covariance of the returns between each pair of stocks is also required. This is estimated using

the sample covariance from the two stocks' return data.

To extend the notation used in the previous sections to two stocks, we simply replace X with Y for the second stock in the pair. The following notation is added:

- (1) Cov_h is the sample covariance determined from using $X_{h,i}$ and $Y_{h,i}$, where $i = 1, 2, \dots, N$.
- (2) $\text{Cov}_a = \frac{\text{Cov}_h}{h}$ is the annualized sample covariance.
- (3) $Z_{i,1}$ are independent standard normal random variables, where $i = 1, 2, \dots, N$.
- (4) $Z_{i,2}$ are standard normal random variables where $i = 1, 2, \dots, N$. They are independent with respect to each other and with respect to the $Z_{i,1}$ random variables.
- (5) ρ is the correlation of the returns of $X_{h,i}$ with $Y_{h,i}$. Since ρ is assumed to remain constant, it does not depend on i (or h).

For the pair of stocks, the *most simple* model for stock returns becomes:

$$X_{h,i} = \mu_X h + \sigma_X \sqrt{h} Z_{i,1}$$

$$Y_{h,i} = \mu_Y h + \sigma_Y \sqrt{h} (\rho Z_{i,1} + \sqrt{1 - \rho^2} Z_{i,2}),$$

where μ_X , μ_Y , σ_X , σ_Y , and ρ are assumed to be constant over the time frame T . The next subsection computes the standard error, $SD[\text{Cov}_a]$, given this model.

4.2 Calculating the standard error

First, determine Cov_a :

$$\begin{aligned} \text{Cov}_a &= \frac{\text{Cov}_h}{h} \\ &= \frac{\sum_{i=1}^N (X_{h,i} - \bar{X}_h)(Y_{h,i} - \bar{Y}_h)}{h(N-1)}. \end{aligned}$$

For large N , this can be re-expressed using the analysis for determining $SD[S_a]$ in the last

section:

$$\begin{aligned} \text{Cov}_a &\sim \frac{\left[\sum_{i=1}^N (\sigma_X \sqrt{h} Z_{i,1}) \right. \\ &\quad \left. \times (\sigma_Y \sqrt{h} (\rho Z_{i,1} + \sqrt{1 - \rho^2} Z_{i,2})) \right]}{hN} \\ &= \frac{\sigma_X \sigma_Y \sum_{i=1}^N (\rho Z_{i,1}^2 + \sqrt{1 - \rho^2} Z_{i,1} Z_{i,2})}{N}. \end{aligned} \quad (3)$$

Taking the variance of both sides of this expression and then noting that $V[Z_{i,1}^2] = 2$ and, since $Z_{i,1}$ and $Z_{i,2}$ are independent, $V[Z_{i,1} Z_{i,2}] = 1$ and $\text{Cov}[Z_{i,2}, Z_{i,1} Z_{i,2}] = 0$, gives

$$\begin{aligned} V[\text{Cov}_a] &\sim \frac{\left[\sigma_X^2 \sigma_Y^2 \sum_{i=1}^N (\rho^2 V[Z_{i,1}^2] \right. \\ &\quad \left. + (1 - \rho^2) V[Z_{i,1} Z_{i,2}] \right. \\ &\quad \left. + 2\rho\sqrt{1 - \rho^2} \right. \\ &\quad \left. \times \text{Cov}[Z_{i,2}, Z_{i,1} Z_{i,2}] \right]}{N^2} \\ &= \frac{\sigma_X^2 \sigma_Y^2 (1 + \rho^2)}{N}. \end{aligned}$$

Finally, taking the square root of both sides and recalling that $T = hN$ yields the standard error:

$$\boxed{SD[\text{Cov}_a] \sim \sigma_X \sigma_Y \sqrt{\frac{h}{T}} \sqrt{1 + \rho^2}}. \quad (4)$$

Note that the $\sqrt{\frac{h}{T}}$ scaling in the expression for $SD[S_a]$ in the previous section also appears here. Should the standard error for the sample variance be desired, it is obtained by simply setting $\rho = 1$ and $\sigma_X \sigma_Y = \sigma^2$ in equation (4).

4.3 Confidence intervals for the covariance

As N grows, the distribution for Cov_a becomes normal, due to applying the central limit theorem to the expression for Cov_a in equation (3). Therefore, the 95% confidence interval for the covariance of the two stocks' returns

becomes $Cov_a \pm 2 \times SD[Cov_a]$ —that is, $Cov_a \pm 2\sigma_X\sigma_Y\sqrt{\frac{h}{T}}\sqrt{1 + \rho^2}$.

Paralleling the analysis for the volatility of a single stock, consider a case where the sample standard deviations of both stocks are 20%, $h = \frac{1}{4 \times 252}$ (sampling four times per trading day), and $T = 1$ (sampling over the time frame of one year). From the analysis for the volatility of a single stock, the volatilities of the two stocks are well approximated by their sample standard deviations, so $\sigma_X \approx \sigma_Y \approx 20\%$. If the sample covariance is 0.03, the sample correlation is 0.75, which may be used to approximate the actual correlation, ρ . Using these values in our confidence interval formula yields 95% confidence that the actual covariance, $\rho\sigma_X\sigma_Y$, is in the interval between 0.02685 and 0.03315. Sampling the returns 40 times, instead of 4 times, per trading day, reduces this interval to being between 0.02900 and 0.03100.

5 Conclusions

The simple calculation in Section 2 clearly shows that historical returns are of little to no practical use in determining the expected return of a stock. At first this suggests looking to determine the expected returns by an alternative method, such as the model of Black and Litterman (1991), but our results also show that there is little to no ability to determine how accurate any alternative method is. That is, our analysis shows that neither historical returns nor future returns are able to provide enough information to reject any hypothesis for the value of the expected return unless it is wildly inaccurate.

Given this, it is important that investment managers using models that employ expected return parameters use them with care, explicitly acknowledging that these parameters cannot be precisely determined. Showing the effect of the expected returns varying over a range of

possibilities can help convey the potential impact of this uncertainty. That said, in some cases like portfolio optimization, this impact can be complicated since the notion of an efficient frontier lacks rigor when expected return values become uncertain.

In contrast, from Section 3, it is clear that the volatility of a stock can be determined reasonably precisely in a relatively short time frame. From Section 4, the same is true for the covariance of two stocks' returns, provided return data is sampled more frequently.

However, these two analyses, as well as the analysis of the expected return in Section 2, assume that the simple return models used in this paper hold. This includes the assumption that the parameters in question, namely the expected return of any stock, the volatility of any stock, and the covariance of any two stocks, all remain constant. As discussed, it is nearly impossible to test if the expected return remains constant, but it is known from return data that the volatility and the covariance of returns, in general, do not remain constant, particularly during bad markets. Therefore, concerns about the validity of the model used here are often larger than concerns about the precision of determining the volatility and the covariance of two stocks' returns from market data.

Endnote

¹ See, for example, https://en.wikipedia.org/wiki/Chi_distribution.

References

- Black, F. and Litterman, R. (1991). "Combining Investor Views with Market Equilibrium," *Journal of Fixed Income* 1(2), 7–18.
- Merton, R. C. (1980). "On Estimating the Expected Return on the Market: An Exploratory Investigation," *Journal of Financial Economics* 8(4), 323–361.

Keywords: Expected return; volatility; covariance; estimation uncertainty; error.