# How Many Realizations Do We Need?

C. V. Deutsch
University of Alberta, Edmonton, Alberta, CANADA (cdeutsch@civil.ualberta.ca)

S. H. Begg
Landmark Graphics Corp., Austin, Texas USA (sbegg@lgc.com)

*A small and seemingly unimportant question comes up toward the end of a geostatistical study of reservoir performance uncertainty:* how many realizations do we need? *A small number of realizations would require less computer resources. A large number should lead to a more reliable assessment of uncertainty. A quantitative method is proposed to calculate the required number of realizations.*

*This first step is to specify the required reliability of the final uncertainty assessment, for example, we require to know the $P_{10}$, $P_{50}$, and $P_{90}$ quantiles within 2%, 4 out of 5 times. The number of realizations can be calculated directly with a requirement expressed in these terms. The exact shape of the response distribution is not needed because reliability is expressed in terms of the cumulative probability. The sampling distribution of the cumulative probability values is Gaussian because the realizations are drawn independently; therefore, we can analytically calculate the required number of realizations.*

*If the number of realizations is fixed because of computational constraints, it is possible to calculate the reliability of the uncertainty predictions. The background of this approach is developed and some examples are presented. Some examples are presented and implementation details are discussed.*

## Introduction and Background

Reservoir management is associated with considerable uncertainty. This uncertainty is due to sparse data and incomplete knowledge of geologic, engineering, and economic factors. Modern decision-making requires an assessment of this uncertainty. A combination of a scenario-based and classical Monte Carlo sampling is used to arrive at probability distributions representing uncertainty in the required performance variables.

The central idea of Monte Carlo sampling is to (1) draw $L$ realizations from a probabilistic model, (2) process the $L$ realizations through some performance calculation, and (3) assemble a histogram of the $L$ responses to represent a distribution of uncertainty in the output(s). Classical Monte Carlo simulation requires the $L$ realizations to be drawn randomly; therefore, they each go into the distribution of uncertainty with equal probability. The critical questions addressed by this short note is: *how many realizations (L) are required?*

The technique developed in this report is general; it is not linked to the particular geostatistical technique used to model facies, porosity, and permeability. Nevertheless, an important discussion in any practical study should focus on how the geostatistical realizations are created since one of the most consequential sources of uncertainty exists in the detailed 3-D distribution of facies and petrophysical properties. Geostatistical techniques are being increasingly used to generate alternative heterogeneous 3-D reservoir models that are consistent with the available data. Decision analysis techniques are being increasingly used to transfer this uncertainty through to reservoir management decision-making.

This note does not address the use of multiple realizations to honor data. There are times when certain realizations would be rejected on the basis of geological expertise or data not used in the geostatistical modeling such as production-related historical observations. We only consider those realizations that meet all basic requirements of reasonableness.

Although a large number of stochastic reservoir models or *realizations* may be available, a small number of realizations are considered in practice. Due to computer limitations, it is only possible to visualize and perform fine-scale full-field flow simulation on a limited number of realizations. Techniques must be applied to reliably choose realizations for more detailed analysis such as flow simulation. A qualitative choice of low, median, and high realizations provides valuable information for reservoir management decision making. Modern decision-making requires a more specific statement of probabilities, e.g., p10, p50, and p90. We would like to identify these limits with as little effort as possible. A companion short note addresses ranking techniques.

We now present a standard methodology for specifying the precision of uncertainty statements, that is, the "uncertainty in our uncertainty statements."

## Precision of Uncertainty Statements

Realizations are drawn for a particular scenario to characterize the uncertainty in reservoir performance variables. There may be multiple scenarios, but we consider one at a time. How many realizations do we need to characterize uncertainty? Figure 1, below, shows 10 realizations. Each realization goes into a histogram of uncertainty with equal probability since geostatistical realizations are equally probably or "equally likely to be drawn."
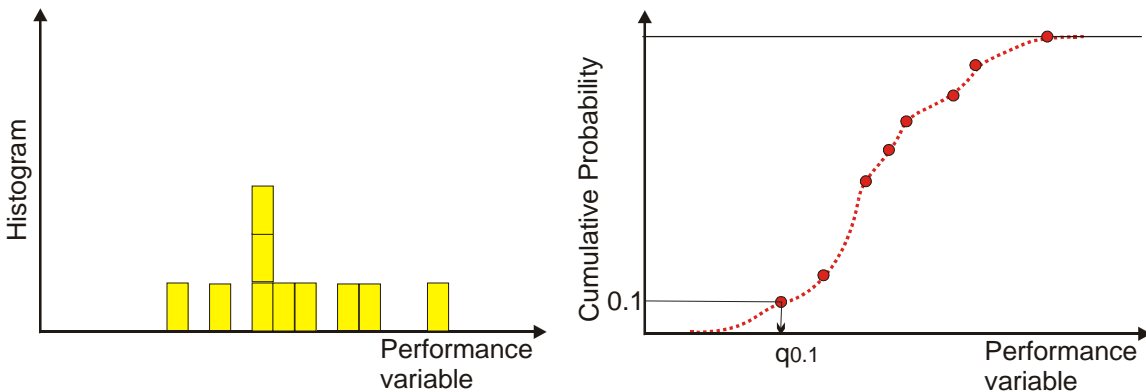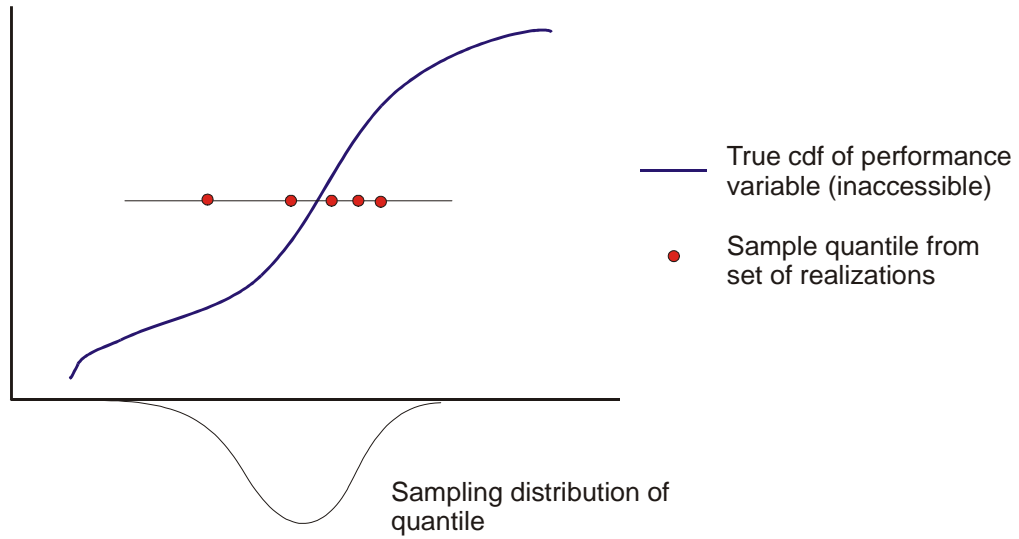


**Figure 1**: 10 realizations of a performance variable shown as a histogram and a cumulative distribution function (cdf). The dashed line on the cdf is an interpretation.

We may want the 0.1, 0.5, and 0.9 quantiles of the distribution of uncertainty. These values can be estimated from the cdf of only 10 realizations, but with great uncertainty or imprecision. It is necessary to specify the required precision in the quantile values, that is, the acceptable "uncertainty in the uncertainty." The 0.1 quantile in Figure 1 would change if 10 different realizations were drawn. There are two different ways of looking at such sampling fluctuations, see Figure 2. The advantages of looking at the performance variable (case A) are that the distribution can be built-up as sets of realizations are drawn and units are easy to interpret. A big advantage of looking at the cumulative probabilities (case B) is that the units are dimensionless. We consider this advantage to be quite important even though we require the cdf of the performance variable (or some reasonable approximation to it). Dynamic determination is considered later.

**A: Sampling distribution of "P"**

True cdf of performance variable (inaccessible)

Sample quantile from set of realizations

Sampling distribution of quantile

**B: Sampling distribution of "F(p)"**

Sampling distribution of F(p)

True cdf of performance variable (inaccessible)

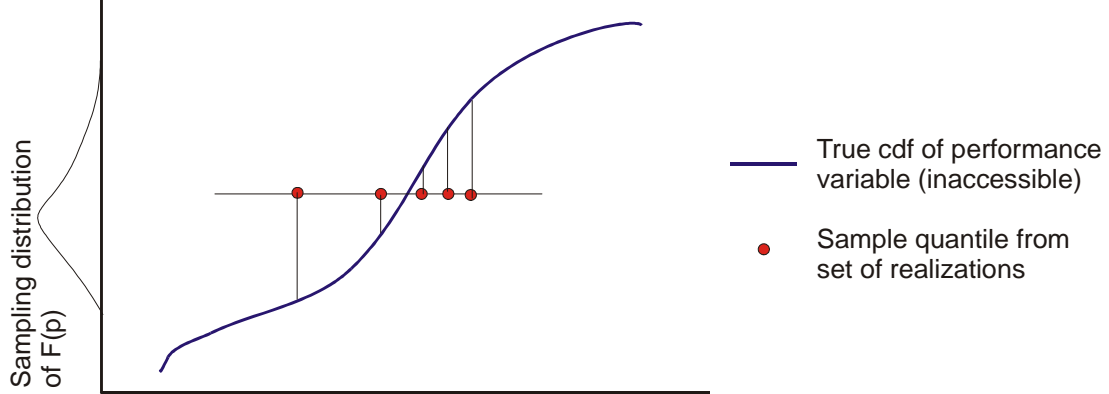Sample quantile from set of realizations

**Figure 2**: Two ways of looking at the sampling distribution of a quantile (A) uncertainty in the quantile, and (B) uncertainty in the real probability associated to the sampled quantile.

Uncertainty assessment is precise if the distribution of cumulative probabilities is narrow. The narrowness could be measured with variance, interquartile range, or some other statistic. A definition consistent with most probabilistic regulatory requirements is considered here. There are two parameters (see Figure 3 for a schematic illustration): (1) $\Delta_F$ = reference difference in probability, which would likely be about 0.01 for reasonable quantiles in the range of 0.1 to 0.9 and smaller, e.g., 0.001 for extreme quantiles in the range of 0.01 to 0.05 and 0.95 to 0.99, and (2) $t_F$ = minimum probability of being within probability $\Delta_F$, which would likely be set to 0.8 or 0.9. The criteria can be expressed as:

> *The estimated F(p)-quantile, p', must be known within a 1% limit*
> *($\Delta_F$ = 0.01) 90% of the time ($t_F$ = 0.9).*

Consider two examples. The first is fairly typical and the second is quite stringent with respect to the quantiles and desired precision:
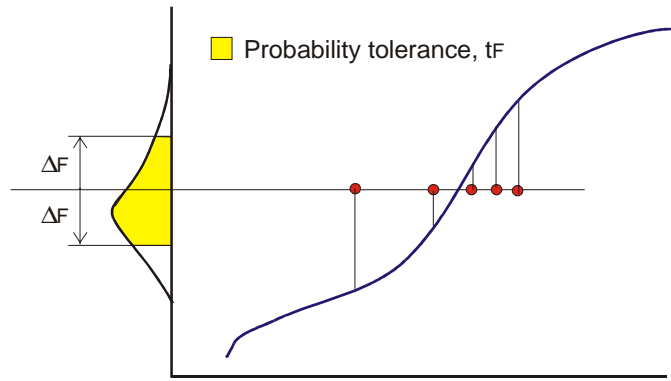
**Figure 3**: Illustration of two parameters ($\Delta_F$ and $t_F$) required to specify required precision of a quantile.

- Require the $P_{10}$, $P_{50}$, and $P_{90}$ of recovery factor with good precision, for example, we require the 0.1, 0.5, and 0.9 quantiles within 1%, 80% of the time ($\Delta_F = 0.01$ and $t_F = 0.8$).

- Require the $P_1$ and $P_{99}$ of net present value with excellent precision, for example, we require the 0.01 and 0.99 quantiles within 0.1%, 95% of the time ($\Delta_F = 0.001$ and $t_F = 0.95$).

The red curve on Figure 4 illustrates how the $t$ parameter changes with the number of realizations (for random or Monte Carlo samples). Recall that $t$ is the proportion of times that the actual probability falls within the specified limit of $\Delta_F$. There would be a family of such curves for different $\Delta_F$ values. The more stringent or the smaller $\Delta_F$, the more realizations that are required, that is, the curve would shift to the right. Successful use of ranking would reduce the number of realizations to achieve a specified $\Delta_F$ tolerance, that is, the curve would shift to the left (see the green curve on Figure 4). We are left to calculate the exact shape and position of the curves illustrated schematically below.
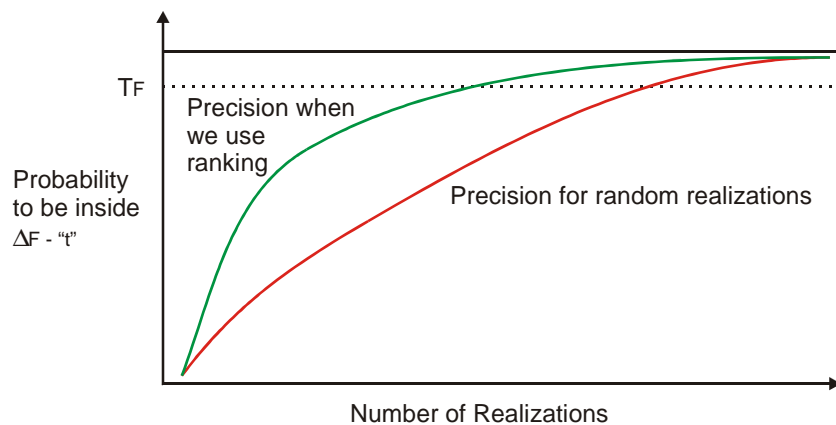


**Figure 4**: Illustration of how precision depends on the number of realizations. The precision will increase more slowly with random realizations than when ranking is used to "target" the realizations toward the quantile.

## Required Number of Realizations

Precision in uncertainty statements is specified by $\Delta_F$ and $t$, that is, a probability tolerance that the quantile should fall within ($\Delta_F$) and the proportion of times the quantile should fall in that tolerance ($t$). This specification of precision is in the units of probability and not the performance variable under consideration. Thus, the analytical and numerical results we derive in this Section are general for all response variables and all applications of Monte Carlo sampling.

Figure 4 presents a functional relationship that appears reasonable, that is, the number of realizations increases as the required precision in our uncertainty specification becomes more stringent. We can establish that relationship by a brute force numerical approach. Consider the following procedure to calculate the value $t$ for a specified $L$, $\Delta_F$, and F values:

1. Draw $L$ realizations or values from a uniform distribution between 0 and 1 (the cdf values on the vertical axes of Figures 2 and 3).

2. Sort those values in ascending order and create a sample cdf. Determine the $F$-quantile of the sample distribution, $F^*$. See if the sample quantile value $F^*$ is within the required tolerance, i.e., $F^* \in [F-\Delta_F, F+\Delta_F[$.

3. Repeat steps 1 and 2 many times (say, $N=10000$) and calculate $t$ as the proportion of times that the sample quantile meets the precision criterion.

This procedure can be repeated for many $L$, $\Delta_F$, and F values to build a family of curves that tell us how many realizations are required for a precision specification. A small program `get_t` was written for this purpose. This program was also customized to report the sampling distribution for the cdf, that is, the distribution of the cdf value (see the distributions on the vertical axes of Figures 2 and 3).

The sampling distributions for the cdf values were found to be normal. This is not surprising since the cdf value is the sum of a large number of values (the indicator transform at the correct threshold) that are independent (the realizations are random) and identically distributed. The central limit theorem tells us that the sampling distribution in this case tends toward a normal distribution as the number increases. The number we are considering ($L$) is very large by central limit theorem standards; therefore, it is expected that the distribution will be normal. Figure 5 shows the sampling distributions for three quantiles and $L=400$. The histograms and probability plot shows that the distribution is normal.

The mean and variance of these distributions can be determined theoretically. The sampled cdf value $F^*$ is the average of an indicator value:

$$F^* = \frac{1}{L}\sum_{l=1}^{L} i(u_l)$$

where the indicator function is 1 if random drawing $u_l$ is less than or equal to the value $F$ and 0 otherwise. The mean or expected value of $F^*$ is the true underlying cdf value, that is, $F$. The variance of $F^*$ is calculated as:

$$\sigma_{F^*}^2 = \frac{1}{L}\sigma_i^2 = \frac{1}{L}\left(E\{i^2\} - (E\{i\})^2\right) = \frac{1}{L}\left(F - F^2\right) = \frac{F(1-F)}{L}$$

These analytical results are verified by the numerical results on Figure 5 and other calculations. The simplicity of these results makes it straightforward to calculate the $t$ statistic for specified $L$, $\Delta_F$, and F values. Recall the definition of the $t$ statistic:

$$t(L, \Delta_F, F) = \text{Prob}\left\{F^* > (F + \Delta_F) \mid L\right\} + \text{Prob}\left\{F^* < (F - \Delta_F) \mid L\right\}$$

The probability distribution of $F^*$ is known to be normal with a mean of $F$ and a variance of (F(1-F))/L; therefore, we can calculate $t$ as a function of the standard normal cumulative distribution function $G(y)$:

$$t(L, \Delta_F, F) = G\left(\frac{\Delta_F}{\sqrt{F(1-F)/L}}\right) - G\left(\frac{-\Delta_F}{\sqrt{F(1-F)/L}}\right)$$

$$= 2 \cdot G\left(\frac{\Delta_F}{\sqrt{F(1-F)/L}}\right) - 1$$

There are numerous "$G(y)$" functions available, e.g., the `gcum` and `ginv` functions in GSLIB. This is an important result. The precision for a given number of realizations can be directly calculated. We can invert this relationship to get an equation that gives us the number of realizations $L$ to achieve a certain precision specification:

$$L = \frac{F(1+F)}{\left(\dfrac{\Delta_F}{G^{-1}((t+1)/2)}\right)^2}$$

This is another important result. The number of realizations for a specified precision can be calculated directly. Figure 6 shows curves of $t$ versus the number of realizations $L$ for $\Delta F = 0.1$, 0.2, 0.3, 0.4, and 0.5. These five curves were also calculated numerically and the results match. These curves relate to 0.5 quantile, which requires the greatest number of realizations ($F(1-F)$ is maximum). Figure 7 shows how the number of realizations changes with quantile; the three curves correspond to $F = 0.1$, 0.25, and 0.5.

## Conclusion

These relations could be implemented in a little calculator-like utility to help people answer questions related to "how many realizations do I need?" and "how good is this $P_{10}$ value?". A little program `numreal` has been written for this purpose. The program is interactive – just run it.

There are no references for this short note; however, workers in statistics must have addressed this problem with similar methods.
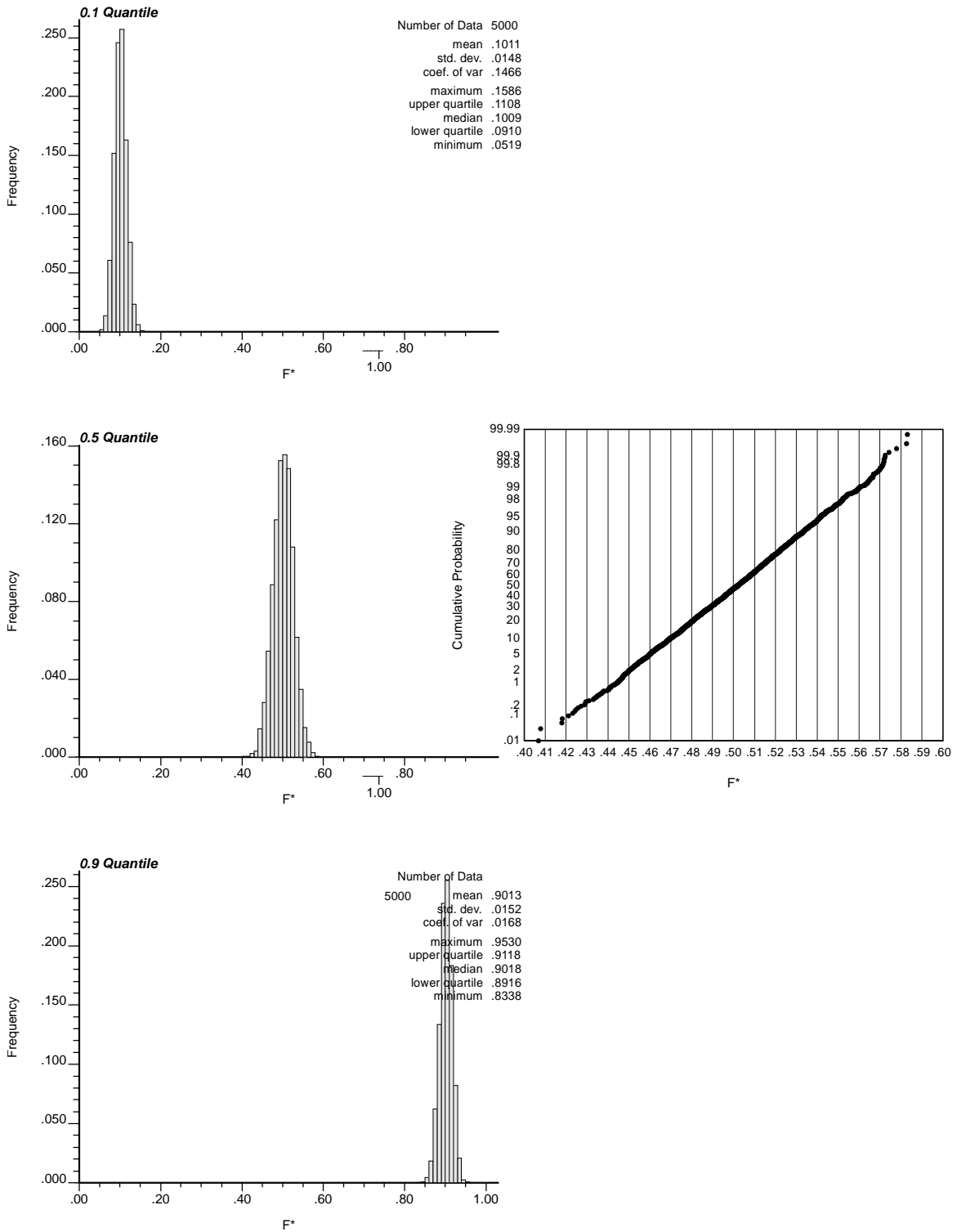
**Figure 5**: Sampling distributions for the 0.1, 0.5, and 0.9 quantile for 400 realizations. The probability plot is shown beside the 0.5 quantile; the straight line indicates a normal distribution.
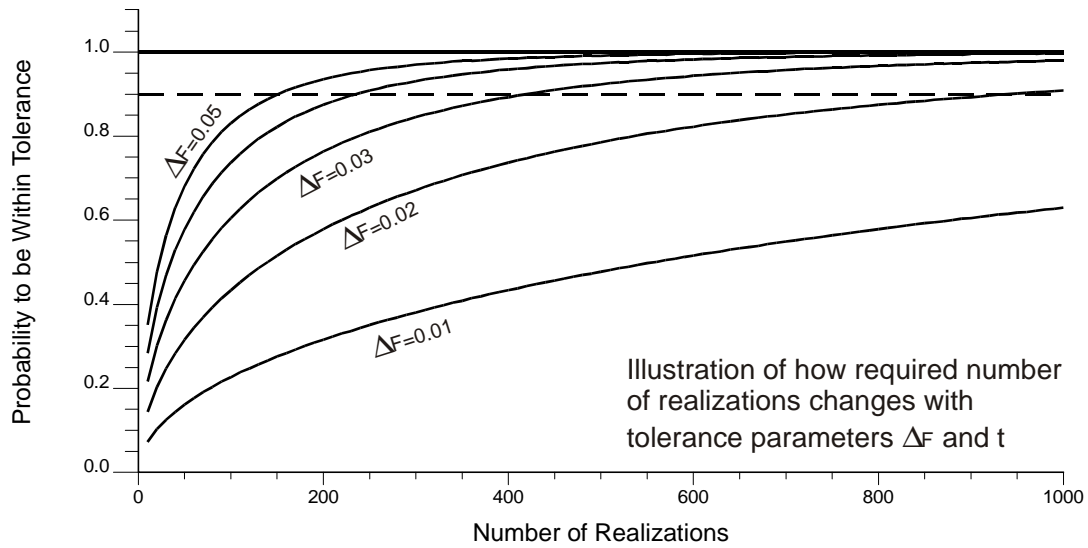
**Figure 6**: Illustration of how the number of realizations changes with tolerance parameters. The four curves correspond to $\Delta F$ = 0.1, 0.2, 0.3, 0.4, and 0.5.
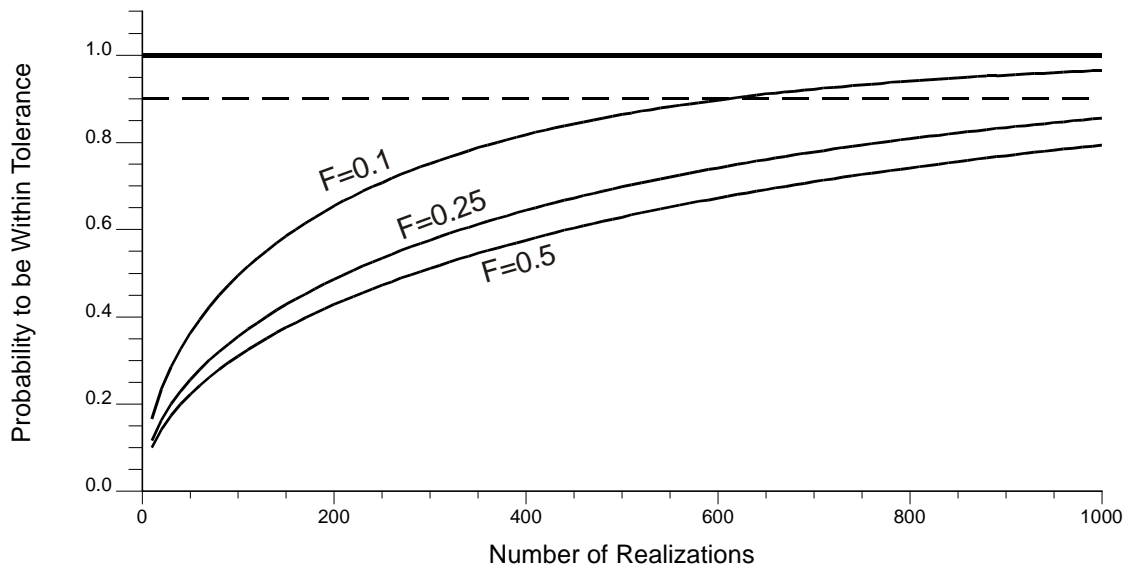


**Figure 7**: Illustration of how the number of realizations changes with quantile. The three curves correspond to F = 0.1, 0.25, and 0.5.