

Why Bayesian Methods Outperform Frequentist Approaches in Project Management

From Historical Databases to Expert-Driven Priors with Bayesian Updating

Murray Cantor, PhD, IBM Distinguished Engineer (Ret.)

March 2026

Executive Summary

For decades, researchers and practitioners have attempted to apply frequentist statistical methods to project management by building databases of project parameters and outcomes. These efforts—spanning defense cost estimation, software engineering benchmarks, transportation megaproject databases, and IT project surveys—have produced useful insights but have consistently run into a fundamental epistemological problem: projects are not drawn from well-defined, homogeneous populations in the way that frequentist inference requires.

The core difficulty is that classical frequentist methods are designed to detect significant differences between mostly similar populations through repeated sampling. Projects, however, are each unique along dozens of dimensions—technology, team composition, political context, regulatory environment, and organizational culture—making it impossible to define reference classes that are simultaneously narrow enough to be meaningful and broad enough to yield adequate sample sizes.

A more natural and epistemologically honest approach uses Bayesian methods: subject matter experts provide three-point estimates (optimistic, most likely, pessimistic) that define triangular prior distributions for task durations and costs. These priors encode contextual judgment that no historical database can capture. As execution proceeds, actual performance data is used to update these priors, producing posterior distributions that narrow over time and converge toward reality. This approach acknowledges high uncertainty in early planning while allowing systematic learning from project-specific data—matching how experienced project managers already think intuitively.

1. The History of Frequentist Methods in Project Management

The application of frequentist statistical methods to project management using empirical databases has a rich, somewhat underappreciated history. What follows traces the major efforts across industries and decades.

1.1 Early Foundations (1950s–1970s)

The RAND Corporation was among the first to systematically collect data on cost and schedule overruns in defense projects. Their work in the 1950s and 1960s laid the groundwork for what would become “reference class forecasting.” PERT (Program Evaluation and Review Technique), developed for the Polaris missile program around 1958, was itself an early attempt to apply probability distributions to task durations, though it used assumed distributions rather than empirical databases.

1.2 Parametric Cost Estimation Databases (1970s–1990s)

The most sustained effort came through parametric cost modeling. Barry Boehm’s COCOMO (Constructive Cost Model), first published in 1981, was built on a database of 63 software projects and used regression analysis to relate project parameters—size, complexity, team capability—to outcomes such as effort and schedule. COCOMO II, published in 2000, expanded this with a larger dataset. Similarly, the PRICE and SEER cost models in defense and aerospace maintained proprietary databases of hundreds of projects, fitting statistical models to predict cost and duration from physical and technical parameters.

1.3 The Standish Group and IT Projects

Starting in 1994, the Standish Group’s CHAOS Reports compiled data on thousands of IT projects, categorizing them as succeeded, challenged, or failed. While their methodology has been criticized—notably by researchers like Robert Glass and Magne Jørgensen—it represented a major attempt to build a frequentist empirical base for IT project outcomes.

1.4 Flyvbjerg and Reference Class Forecasting

Perhaps the most rigorous and influential effort came from Bent Flyvbjerg, starting with his work in the early 2000s. His databases covered transportation infrastructure projects—roads, bridges, rail, tunnels—and later expanded to other megaproject types. His key 2002 paper with Holm and Buhl analyzed 258 transportation projects across 20 countries spanning several decades. The core finding was systematic cost overruns with identifiable statistical distributions.

Flyvbjerg explicitly advocated “reference class forecasting,” where you place a new project into a class of completed similar projects and use the empirical distribution of outcomes as your forecast. This was adopted as official policy by the UK Treasury in 2003 and later by the Danish

and Swiss governments. His later book, *How Big Things Get Done* (2023), drew on a database of over 16,000 projects across multiple sectors.

1.5 NASA and Aerospace

NASA maintained several databases, including the NASA Cost Estimating Handbook data and the Air Force Cost Analysis Agency databases. These supported parametric models like the NASA/Air Force Cost Model (NAFCOM), which used regression on historical mission data to estimate costs for new missions based on weight, technology readiness, and other parameters.

1.6 Software Engineering Empirical Databases

The International Software Benchmarking Standards Group (ISBSG) has maintained a database of thousands of software projects since the mid-1990s, collecting function point counts, effort, defects, and other metrics. The idea was explicitly frequentist: given a new project's parameters, look up the empirical distribution of outcomes for similar past projects. The Experience Factory concept from Victor Basili at the University of Maryland similarly promoted the systematic collection and reuse of project data.

1.7 Construction Industry

In construction, databases like those maintained by RSMeans (now Gordian) have provided empirical cost data for decades. The UK's Building Cost Information Service (BCIS) has similarly compiled project cost data to support statistical estimation.

2. The Fundamental Problem with Frequentist Approaches

Classical frequentist inference is built around the idea of repeated sampling from a well-defined population. The machinery of confidence intervals, p-values, and hypothesis tests all assume you can meaningfully say “this observation is drawn from the same generating process as these other observations.” That works beautifully when comparing drug A to placebo across thousands of patients, or testing whether two manufacturing lines produce different defect rates.

Projects, however, resist this framing in a fundamental way. Each major project is arguably unique along dozens of dimensions simultaneously—technology, team, political context, regulatory environment, stakeholder dynamics, organizational culture. When Flyvbjerg groups “rail projects” together, he is making a judgment call that these projects share enough of a common data-generating process to justify pooling. But a high-speed rail line in Norway and a metro extension in Mumbai may share little beyond the presence of tracks.

The core problem is not just that you cannot form populations—it is that you can never be confident your reference class is homogeneous enough for the frequentist machinery to be valid. You are always caught in an inescapable tension: make the reference class broad enough to get a usable sample size, and you have pooled dissimilar things; make it narrow enough to be genuinely comparable, and you have five projects, which tells you almost nothing statistically.

The parametric cost models (COCOMO, NAFCOM) tried to sidestep the population problem through regression—treating project parameters as predictors rather than defining discrete populations. But that merely relocates the problem to the assumption that the functional relationship between parameters and outcomes is stable across the database, which is also questionable.

2.1 Persistent Challenges

These efforts have faced persistent difficulties across all domains. Projects are arguably non-identical, making “reference classes” debatable. Selection bias is rampant—failed projects are underreported. The dimensionality of project characteristics is high relative to available sample sizes. And organizational context—culture, governance, political dynamics—is extremely difficult to parameterize. Even Flyvbjerg’s own work has faced methodological scrutiny, including questions about how reference classes are defined and whether his overrun distributions are robust to different sampling choices.

3. The Bayesian Alternative: Expert Priors with Bayesian Updating

A more natural approach begins with subject matter experts. The three-point estimate—optimistic, most likely, pessimistic—is essentially an elicited prior distribution. When you fit a triangular or PERT-Beta distribution to those three points, you are encoding one person’s informed judgment about the range of plausible outcomes for a specific task or cost element. This is a fundamentally Bayesian act, even when practitioners do not use that language. You are not claiming this task was drawn from a population of similar tasks. You are saying “given everything this expert knows about this particular task in this particular context, here is their uncertainty.”

3.1 Why This Fits Projects Better

The subject matter expert bringing their judgment to bear is doing something no historical database can do—they are integrating the idiosyncratic features of this project. They know the team is inexperienced, or that the client changes scope frequently, or that the technology is unproven but the vendor seems competent. That contextual richness is precisely what gets lost when you try to pool projects into reference classes.

The triangular distribution is popular for good reason. It is intuitive for experts to parameterize, it accommodates asymmetric uncertainty (most projects have longer right tails—things go wrong more dramatically than they go right), and it requires no assumption that the task belongs to a well-defined statistical population.

3.2 The Update Cycle

Once execution begins, you start collecting actuals—real expenditures, real durations, real earned value. At this point you can update the prior in a principled way.

Consider an example: a subject matter expert estimates a work package will take 40 to 80 days, most likely 55. You are now 30 days in and 40 percent complete. The burn rate data is telling you something the original estimate did not know. In a Bayesian framework, you combine the prior (the expert’s triangular distribution) with the likelihood (what the tracking data implies about the true parameter) to get a posterior distribution that is typically narrower and shifted toward what the actuals suggest.

In practice this often happens through Monte Carlo simulation. Tools like Primavera Risk Analysis or Risky Project run thousands of iterations, sampling from the prior distributions on remaining work, but those priors can be tightened or shifted as actuals accumulate. The key insight is that the prior does not have to be right—it just has to be a reasonable starting point that the data can correct.

3.3 The Practical Workflow

This gives you a natural project lifecycle for uncertainty management. In early planning, when you have no actuals, you are entirely dependent on expert judgment. The distributions are wide and the forecasts are honestly uncertain. As you move through execution, each reporting period brings new data, and the forecasts should narrow and become more reliable. This is exactly the behavior you want—high acknowledged uncertainty early, converging toward reality as you learn.

Compare this to the frequentist alternative: you would need to find a database of completed similar work packages, verify the reference class is appropriate, and apply the historical distribution to your task. But “similar” is doing enormous work in that sentence. The Bayesian approach acknowledges that the expert’s brain is the best available integrator of all the relevant contextual factors, and then lets the data discipline that judgment over time.

3.4 Where Historical Data Still Helps

This does not mean empirical databases are useless—they play a valuable role in calibrating and challenging expert estimates. If an expert says a task will take 20 to 30 days and the historical base rate for similar tasks shows a median of 45 days with heavy right skew, that is a signal the expert may be anchored too optimistically. Flyvbjerg’s “outside view” data is genuinely useful as a sanity check on the “inside view” of project-specific estimation. The point is that the historical data informs the prior rather than replacing expert judgment entirely.

4. Conclusion

The history of applying frequentist methods to project management is largely a history of people trying to force project data into a statistical framework whose assumptions do not hold. Projects are not repeated draws from homogeneous populations. Reference classes are always debatable. Sample sizes within meaningful categories are always too small.

The Bayesian alternative—eliciting expert priors through three-point estimates and updating them with actuals as execution proceeds—is more epistemologically honest and more practically useful. It acknowledges uncertainty where it exists, incorporates contextual judgment that databases cannot capture, and provides a principled mechanism for learning from project-specific data over time.

Despite being the more natural framework for project estimation, this approach remains underused in many organizations. Project managers often learn PERT scheduling without understanding the probabilistic foundations. Organizations frequently demand single-point estimates rather than distributions. And the update discipline—actually revising forecasts as actuals come in rather than just reporting variance—requires a maturity in project controls that many teams have not yet achieved.

The irony is that most experienced project managers already think in something like this framework intuitively. They ask experts for ranges, they revise expectations as work progresses, and they worry more about estimates where the expert seemed uncertain. The Bayesian formalization simply makes that natural reasoning rigorous, transparent, and auditable.

References

- Basili, V.R., Caldiera, G., and Rombach, H.D. (1994). "The Experience Factory." In *Encyclopedia of Software Engineering*, vol. 1, pp. 469–476. John Wiley & Sons.
- Boehm, B.W. (1981). *Software Engineering Economics*. Prentice-Hall, Englewood Cliffs, NJ.
- Boehm, B.W., Abts, C., Brown, A.W., Chulani, S., Clark, B.K., Horowitz, E., Madachy, R., Reifer, D.J., and Steece, B. (2000). *Software Cost Estimation with COCOMO II*. Prentice-Hall.
- Flyvbjerg, B., Holm, M.S., and Buhl, S. (2002). "Underestimating Costs in Public Works Projects: Error or Lie?" *Journal of the American Planning Association*, 68(3), pp. 279–295.
- Flyvbjerg, B. (2006). "From Nobel Prize to Project Management: Getting Risks Right." *Project Management Journal*, 37(3), pp. 5–15.
- Flyvbjerg, B. and Gardner, D. (2023). *How Big Things Get Done: The Surprising Factors That Determine the Fate of Every Project, from Home Renovations to Space Exploration and Everything in Between*. Crown Currency.
- Glass, R.L. (2006). "The Standish Report: Does It Really Describe a Software Crisis?" *Communications of the ACM*, 49(8), pp. 15–16.
- HM Treasury. (2003). *The Green Book: Appraisal and Evaluation in Central Government*. TSO, London.
- International Software Benchmarking Standards Group. (2018). *ISBSG Software Development and Enhancement Repository*, Release 2018. ISBSG Ltd.
- Jørgensen, M. and Moløkken-Østfold, K. (2006). "How Large Are Software Cost Overruns? A Review of the 1994 CHAOS Report." *Information and Software Technology*, 48(4), pp. 297–301.
- Kahneman, D. and Tversky, A. (1979). "Prospect Theory: An Analysis of Decision Under Risk." *Econometrica*, 47(2), pp. 263–291.
- Malcolm, D.G., Roseboom, J.H., Clark, C.E., and Fazar, W. (1959). "Application of a Technique for Research and Development Program Evaluation." *Operations Research*, 7(5), pp. 646–669.
- NASA Cost Estimating Handbook. (2015). NASA Office of the Chief Financial Officer, Cost Analysis Division. Washington, DC.
- Project Management Institute. (2017). *A Guide to the Project Management Body of Knowledge (PMBOK Guide)*, 6th edition. PMI, Newtown Square, PA.
- RAND Corporation. (1993). *A Review of the RAND Corporation's Research on Military Cost Analysis*. RAND, Santa Monica, CA.
- The Standish Group. (1994). *The CHAOS Report*. The Standish Group International.
- Vose, D. (2008). *Risk Analysis: A Quantitative Guide*, 3rd edition. John Wiley & Sons.