

Becoming Bayesian

Making Decisions with Bayes Factors in Flight Test

Brig Gen Douglas “Beaker” Wickert, PhD¹

Air Force Materiel Command

*The fundamental cause of trouble in the modern world is that
the stupid are cocksure while the intelligent are full of doubt*

Bertrand Russell [1]

Decision-making in flight test is done under a persistent state of uncertainty. The risk management frameworks used in flight test are typically probability-based, but assessing probability is largely subjective due to the epistemic uncertainty inherent in flight test. Bayesian reasoning, a method that continuously updates our beliefs about the world in light of new evidence, is widely used in scientific fields. This paper describes an approach using *Bayes Factors* to support our technical and safety risk assessments in test. The use of Bayes Factors provides a systematic structural framework with a common lexicon, leveraging the collective wisdom of the crowd during test and safety planning, to enable meaningful conversations about risk. Multiple examples of the practical use of Bayes Factors are provided. Lessons learned and recommendations for flight test risk management are offered.

Keywords: Bayes Factor, Risk, Risk Management, Flight Test, Uncertainty

Nomenclature

BF	Bayes Factor (likelihood ratio)
$\Theta^{(0)}$	prior odds ratio
$\Theta^{(1)}$	posterior odds ratio
H	hypothesis
$\neg H, H^c$	negation or complementary (opposite) hypothesis
D	data
$P(H D)$	conditional probability of hypothesis given the data

¹ Associate Fellow SETP

I. Uncertainty & Flight Test

Uncertainty is unavoidable and ubiquitous. This is more true in flight test than in many other professions. To reduce uncertainty about a system under test, flight testing deliberately probes the unknown, “pushing the edge of the envelope,” identifying deviations in system performance from the intended design. Test also reveals unexpected behavior, model flaws, system deficiencies, poor human factor designs, and inaccurate design assumptions. Due to that inherent uncertainty, flight test is risky and potentially dangerous. A century of flight test has led to the development of a robust approach to risk management and a systematic, deliberate way of approaching the unknown—*Ad Inexplorata*² was adopted in 1953 as the motto of the Air Force Flight Test Center.³ The “X” logo of the Society of Experimental Test Pilots, established in 1955, stands for the *unknown*.

Flight testers are professional risk managers. Over the last several years, a new approach to considering uncertainty and communicating about risk has emerged at Edwards. Built around the Bayesian reasoning techniques that are now common in scientific data analysis, the approach has proven to be useful for thinking about risk and communicating with a common framework.

In many ways, this paper is the next step in a career spent thinking about risk management and uncertainty. In a 2018 SETP paper, the author attempted to describe the attributes of a risk-aware culture and techniques for cultivating **Risk Awareness** [2]. Risk Awareness was defined in an analogy to Situational Awareness: the *perception of the elements of uncertainty and the potential, projected outcomes resulting from uncertainty*.⁴ In contrast to the MIL-STD-882E definition of risk, “the combination of the severity of a mishap and the probability of that mishap occurring” [4], the author deliberately avoided considering ‘probability-of-occurrence’ in the definition of Risk Awareness because of the inherent unreliability of probabilities. This is particularly true given the epistemic uncertainties present in flight test. Since 2018, the author has ‘*Become Bayesian*’ and has embraced Bayes Factors as a way of dealing with the inherent difficulty in assessing probabilities. This paper describes the approach.

In his reflection on a “whole career working on investigations aimed at reducing uncertainty about what is happening, what might happen, and even the reasons why things happen,” the eminent statistician Sir David Spiegelhalter, Professor Emeritus at Oxford and past President of the Royal Statistical Society, writes that: “once we accept a personal, subjective view of probability and uncertainty, we are led naturally to *Bayesian analysis*, in which we use the theory of probability to revise our beliefs in the light of new evidence” [5].

In simple terms, a Bayesian approach is one in which we continuously update our beliefs about the world in the light of new evidence. A little more formally and in the terms used in Bayesian inference:

² “Toward the Unexplored”

³ The Air Force Flight Test Center (AFFTC) was redesignated as the Air Force Test Center (AFTC) in 2012 as part of a major reorganization of Air Force Materiel Command; the reorganization brought AFFTC, the Arnold Engineering Development Center, and the test organizations in the Air Armament Center together into a single Center.

⁴ Situational awareness is defined as the elements in the environment within a volume of time and space, the comprehension of their meaning, and the projection of their status in the near future [3].

$$\begin{pmatrix} \text{Prior} \\ \text{Beliefs} \end{pmatrix} \times \begin{pmatrix} \text{Likelihood of} \\ \text{new Evidence} \end{pmatrix} \Rightarrow \begin{pmatrix} \text{Posterior} \\ \text{Beliefs} \end{pmatrix} \quad (1)$$

This paper describes an approach of using Bayes Factors (an assessment of the likelihood of observing a set of evidence) to update our prior judgment on whether a proposed test has reached the threshold for an acceptable level of risk (ALR). The evidence may be engineering predictions, test data, unexpected test events, test team training and readiness, concerns of programmatic drift, *etc.* The use of Bayes Factors also provides a structural framework with a common lexicon during test and safety planning to facilitate meaningful conversations about risk. This can help leverage the collective “wisdom of the crowd” during test planning and execution.

Before describing the formal technique of using Bayes Factors in Section III, we begin with a description of probability, judgment, and decision-making under uncertainty (Section II). Section II may be skipped without any loss in understanding the application of Bayes Factors, but the description of a century-plus-long intellectual journey through uncertainty and decision-making contributes to understanding the theoretical foundation of our recommendations for making good decisions. After describing the Bayes Factor approach, the paper includes eight applications of the use of Bayes Factors using both test and non-test examples (Section IV). Each example illustrates a different attribute and lesson learned from ‘*Becoming Bayesian*,’ which are discussed and summarized in Section V. Two appendices provide the mathematical basis for the derivation of Bayes Factors from Bayes’ Theorem (Appendix A) and of the use of the Beta and Gamma distributions in Bayesian reasoning (Appendix B).

II. Risk & Decision-Making Under Uncertainty

Decision-making under uncertainty is a fundamental challenge encountered in many human endeavors, from medical diagnosis and financial investment to military strategy and engineering design. A rich body of literature exists on the subject [6, 7, 8]. Uncertainty is particularly ubiquitous and persistent in flight test and is often a topic at SETP and SFTE symposia.

The author’s preferred definition of uncertainty—*uncertainty* is the *conscious awareness of ignorance*—comes from Spiegelhalter’s excellent book on the topic [5]. Our ignorance—our lack of knowledge of “what we know” or “what we can know”—stems from two primary factors: the theoretical limits of knowledge (*epistemic*⁵ uncertainty) and the inherent variability or randomness in some systems (*aleatory*⁶ uncertainty). Ignorance is our natural state, and as James Clerk Maxwell noted, “Thoroughly conscious ignorance is the prelude to every real advance in science” [9].

Risk is a fundamental consequence of uncertainty. The International Organization for Standardization (ISO) defines risk as “the effect of uncertainty on objectives” [10], while the Air Force defines it as “probability and severity of loss or adverse impact from exposure to various hazards” [11]. The definition used by the Air Force Test Center parallels the Air Force definition: test risk is the “combination of the severity of the mishap and the probability that the mishap will occur”

⁵ From the Greek *episteme*, meaning “knowledge” or “understanding.”

⁶ From the Latin root *alea*, meaning “dice” or “game of chance;” in Latin, an *aleator* was “a dice player.”



Figure 1: Risk Management Framework and Associated Analysis Tools

[12].

All definitions of risk effectively have three elements in common: 1) scenarios; 2) consequences; and 3) likelihoods. ISO's definition of risk management is the "coordinated activities to direct and control an organization with regard to risk." Similarly, the Air Force defines risk management as a "decision-making process to systematically evaluate possible courses of action, identify risks and benefits, and determine the best course of action" [13]. Risk management is equivalent to answering three questions related to the three elements in Figure 1:

1. What can go wrong? (*generally an undesired/unexpected event*)
2. What are the resulting consequences? (*negative outcomes from the unexpected event*)
3. What is our expectation of the likelihood? (*probability of occurrence*)

The author has long been uncomfortable with the third question. This is primarily due to the difficulty in making probability assessments with any confidence (Section II.A). The emphasis in the Risk Awareness framework [2] on uncertainty was an attempt to avoid the question of likelihood. However, as discussed in Section II.C, we implicitly make probability assessments every time we make decisions under uncertainty. Thus, probability is an unavoidable and integral part of risk management.

Within the three-element conceptual model of risk in Figure 1, Risk Awareness supports, guides, and informs the scenario analysis and planning process. Test Hazard Analysis (THA), Failure Mode and Effects Analysis (FMEA), System-Theoretic Process Analysis (STPA) [14], and other system engineering tools are useful in defining the scope and impact of consequences.⁷ Many tools exist to support the third element: Monte Carlo methods, Probability Risk Assessment [15], and expert opinions are examples of approaches to assess likelihoods. In this paper, we offer Bayes Factors as another. Our adoption and approach to using Bayes Factors to inform risk management is motivated by two axioms that serve as guiding principles:

1. Objective probability does not exist (but it is useful to act as if it does) (Section II.A)
2. Harness the *Wisdom of the Crowd* (the fact that aggregate judgments can often outperform individual assessments by experts) (Section II.B)

⁷ Often, STPA and THA will start with consequences and identify scenarios and hazards that could result in the consequences; in the case of STPA, this is done with the goal of "engineering away" the consequence.

A. Probability does not exist

Probability estimates are essential to risk management and decision-making under uncertainty. Given the unavoidable need to assess probability, it is worth noting that we do not even have a good way to define it. Bertrand Russell, whom we already quoted in the epigraph, remarked in a 1929 lecture that “Probability is the most important concept in modern science, especially as nobody has the slightest notion what it means” [16].

There is a long and rich history of debate on the nature of probability. The major interpretations and attempts to define probability since Pascal and Fermat first corresponded⁸ on the topic in the late 1600s have included: *classical* probability (Laplace [17], who built on Bernoulli’s law of large numbers [18]), *logical/evidential* probability (Keynes [19], Carnap [20]), *subjective* probability (Bayes [21], de Finetti [22]), *frequency* interpretations (von Mises [23], Fisher [24], Neyman [25]), *propensity* interpretations (Popper [26]), and *axiomatic* definitions (Kolmogorov [27]). Rehashing the metaphysical debate is not our current purpose.⁹ But awareness of the difficulty to even define *probability* means that we should expect difficulty in assessing it.

The author finds the most useful statement on probability to be consistent with Italian mathematician Bruno de Finetti’s perspective that “Probability does not exist” [22]. By this, de Finetti meant that probability is not an inherent property of the world, but rather a degree of belief held by an individual. This view was controversial and opposed by both *frequentists* like Fisher, Neyman, and Pearson, who saw probability as long-run frequency of expectations, and *objectivists* like Keynes, who saw probability as the objective logical relation between propositions [19]. Ramsey, whom we will meet again below, was sharply critical of Keynes’ view and also proposed a subjective interpretation in which probability reflects an individual’s coherent degrees of belief [29]. This disagreement in the nature of probability would be entwined with the bitter debate that raged for much of the 20th century between *frequentists* and *Bayesians*.

Today, the philosophical and academic debate is over for all practical purposes, with most practitioners¹⁰ adopting a Bayesian understanding that rests on de Finetti’s subjective interpretation of probability [30, 5, 31, 32, 33, 34]. Probability is best understood as our subjective relationship with what we know about the world, and Bayesian reasoning is the process by which we update our subjective understanding. **Probability is a direct measure of our uncertainty.** If I flip a coin and ask you about *your probability* of the heads/tails outcome, your expectation is unaffected by whether I look at the result after the flip. ‘*My probability*’ given knowledge of the outcome is obviously different from ‘*your probability*’ of heads/tails without my knowledge of the outcome.¹¹ Our different measures of uncertainty result in different assessments of the probability.

Estimating probabilities in the real, ambiguous world is notoriously challenging. For example, what is the probability of interest rates being below 2.5% a year from now? Ten years

⁸ A collaboration sparked by a problem presented to Pascal by Parisian gambler and writer, the Chevalier de Méré, who wanted to know how to divide the stakes of a prematurely interrupted game fairly.

⁹ There are many excellent summaries of the debate; see [28] for example.

¹⁰ In addition to being standard in data analysis, Bayesian techniques are widely used in AI and machine learning; large-language models are typically built on Bayesian foundations.

¹¹ This example also illustrates the distinction between aleatory and epistemic uncertainty: before the flip, there is aleatory uncertainty of the outcome; after the flip, my knowledge of the outcome does not affect your residual epistemic uncertainty of the outcome.

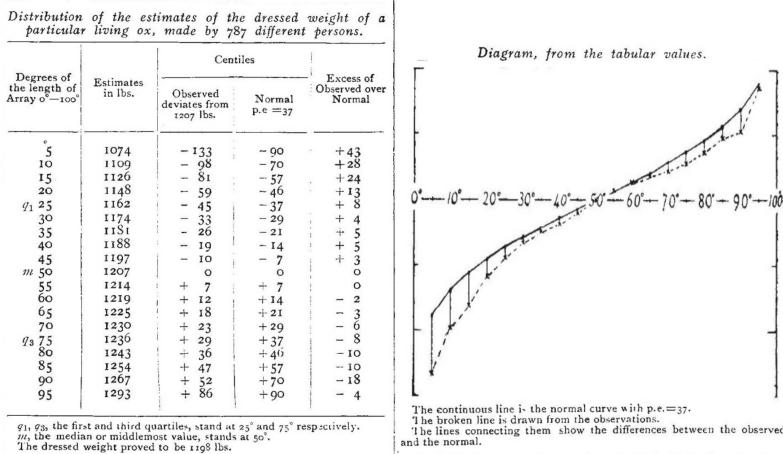


Figure 2: Figures from Sir Francis Galton's 1907 Nature article noting the wisdom of crowds

from now? What is the probability that Chairman Xi decides to invade Taiwan? What is the probability that the Atlantic Meridional Overturning Circulation (AMOC)¹² collapses in the next twenty years? In the next fifty years?

Engineering predictions can also be notoriously unreliable. For example, reliability calculations by McDonnell Douglas had estimated that the probability of a DC-10 experiencing a simultaneous loss of engine thrust and asymmetric leading edge slats to be less than one-in-a-billion (10^{-9}) [35]. But this identical combination of failures occurred four times in four years and was the root cause of the 1979 crash at Chicago O'Hare, still the deadliest single airline accident on US soil. Similarly, various probability estimates by NASA engineers for loss of the Space Shuttle ranged from 1-in-10 to 1-in-100,000, a discrepancy spanning five orders of magnitude [36, 37]. More recently, the flight testing of a new platform at Edwards AFB experienced a 10^{-9} anomaly during the first 10 hours of test operations.¹³

If we accept probability as a measure of our uncertainty, then probability does not objectively exist. Yet we need probability judgments for decision-making under uncertainty, so it is useful to act as if they do actually exist. Bayes Factors offer one mechanism for expressing our uncertainty and communicating our judgment of probability. And, as we shall see, Bayes Factors provide a framework that can leverage the wisdom of the crowd, with the goal of improving our collective assessments and decision-making.

B. Wisdom of the crowd

The “wisdom of the crowd” is the phenomenon where aggregated judgments from diverse groups outperform individual assessments by experts. The principle was first quantitatively observed by Sir Francis Galton at a 1906 country fair in Plymouth, England [38]. When Galton

¹² AMOC is the “conveyor belt” moving warm, salty water from the tropics northward and colder, deeper water southward; the circulation is a crucial component of the global climate system.

¹³ If the aircraft flew for 8 hours a day, 365 days a year, the mean time for the first such expected failure should be 100 years.

analyzed nearly 800 guesses from fairgoers attempting to estimate the weight of an ox, he discovered that while individual estimates varied widely, the median of all guesses was remarkably accurate—within one percent of the actual weight of 1,198 pounds (Figure 2). This result has been repeatedly validated: collective intelligence often emerges from the aggregation of diverse, independent judgments, even when individual participants lack specific expertise.

A century after Galton, Philip Tetlock’s Good Judgment Project provided rigorous empirical validation of the principle of the wisdom of the crowd through large-scale forecasting tournaments in which teams of ordinary citizens, when properly structured and trained in probabilistic reasoning, outperformed intelligence analysts with access to classified information [39]. Tetlock identifies the conditions for effective crowd wisdom: diversity of perspectives, independence of initial judgments, and an aggregation mechanism. Through the Good Judgment Project, Tetlock and journalist Dan Gardner claim to have also identified ‘superforecasters’ who consistently outperform both crowds and experts [40]. Fittingly to our present purpose, one of the common attributes of superforecasters’ success was “being Bayesian,” *i.e.*, regularly updating their beliefs in response to new information.¹⁴

In *The Wisdom of Crowds* James Surowiecki describes four conditions necessary for crowd wisdom to emerge: diversity of opinion among participants, independence of individual judgments (avoiding cascading influence), decentralization that allows people to draw on local knowledge, and an effective aggregation mechanism to combine individual inputs into collective decisions [42]. These are also worthwhile guides for our flight test safety review processes.

In the author’s experience, the flight test community does a good job of considering and synthesizing judgments from diverse teams—as pilots, engineers, maintainers, flight doctors, and independent safety reviewers typically participate in comprehensive safety reviews. If improvement is needed, it is in having a framework and common lexicon for meaningful conversations about risk. Without a common lexicon, it is difficult for experts from different fields to communicate and compare opinions on relative likelihoods. This is one of our motivations for encouraging a widespread adoption of Bayes Factors.

C. Making Decisions Under Uncertainty

Managing risk is essentially the same as making decisions under uncertainty, and in making decisions under uncertainty, we are either implicitly or explicitly making probability judgments. The

¹⁴ Other attributes include good qualities that all testers wanting to improve decision-making should aspire to:

- *Aggregation*: systematically using multiple sources of information
- *Openness*: openness to new knowledge and actively seeking contrary views and perspectives
- *Meta-cognition*: insight into their own thinking and awareness of cognitive biases
- *Decomposition*: breaking complex problems into smaller components more tractable to analysis
- *Humility*: acknowledging uncertainty, expressing predictions probabilistically rather than definitively, and readily admitting errors
- *More ‘Fox’ than ‘Hedgehog’*: from Isaiah Berlin’s 1953 essay: “the fox knows many things, but the hedgehog knows one big thing” [41]; describes a fundamental difference in thinking styles—*foxes* are generalists with many diverse, conceptual models; *hedgehogs* are specialists who relate everything to a single idea or system

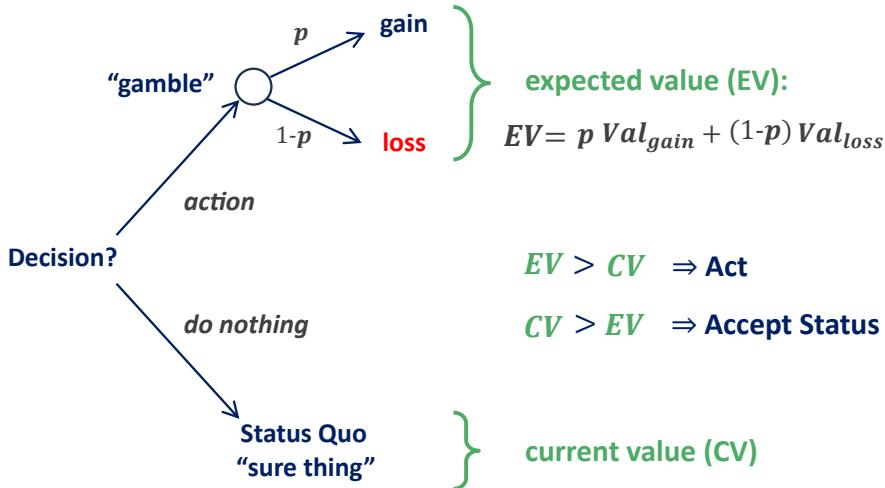


Figure 3: Ramsey’s Decision Theory Calculus

psychology of decision-making is a deeply researched field. An excellent recent paper [43] surveys the field from the Greeks through 21st-century thinkers, including Kahneman and Gigerenzer. Risk Awareness [2] described in detail the complementary perspectives of decision-making under uncertainty offered by Kahneman [44] and Gigerenzer [45]. Here, we provide a brief overview of the major 20th-century intellectual perspectives on risk and decision-making.

Knight (1921):¹⁵ first articulated the distinction between *risk*—situations where the outcomes are unknown but governed by probability distributions—and *uncertainty*—situations in which outcomes and probability models are unknown [46]. In modern terms, we refer to aleatory and epistemic uncertainty. Risk Awareness [2] argued that different cognitive and risk management tools were needed for the different risk and uncertainty domains. Stirling made a similar argument for policy considerations in an influential letter in Nature [47].

Ramsey (1926):¹⁶ Ramsey’s decision calculus formalized the idea that rational agents choose actions that maximize their expected utility, where utility represents the subjective value of outcomes and probabilities represent degrees of belief (Figure 3) [29, 48]. Ramsey’s decision theory provides the mathematical foundation for treating uncertainty as a quantifiable factor in decision-making and establishes a theoretical basis for a Bayes Factor approach to risk management.

Simon (1947):¹⁷ introduced *bounded rationality* and argued that human decision-makers, faced with cognitive limitations, incomplete information, and time constraints, engage in *satisficing* behavior—searching through alternatives until finding one that meets their minimum acceptable criteria (Figure 4) [49, 8]. Decision makers often stop searching when they find an option that *satisfices* some objective criteria. There is an implicit, subconscious probability assessment of the likely outcome when one stops searching through the space of bounded rationality.

¹⁵ Frank Knight (1885–1972): American economist and one of the founders of the Chicago School.

¹⁶ Frank P. Ramsey (1903–1930): a brilliant British mathematician, philosopher, and economist whose tragically brief career left foundational contributions to decision theory that remain central to modern risk analysis.

¹⁷ Herbert Simon (1916–2001): a Nobel Prize-winning American economist and cognitive scientist.

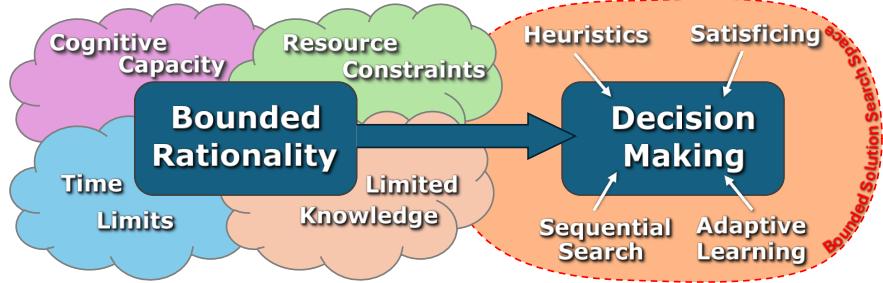


Figure 4: Simon’s Bounded Rationality

Kahneman and Tversky (1974):^{18,19} in the highly influential paper, “Judgment Under Uncertainty: Heuristics and Biases,” they highlighted three types of heuristics by which probabilities are assessed and judgments made: availability, representativeness, and anchoring and adjustment [6]. Kahneman spent his career systematically revealing the flaws in human reasoning and introduced the dual-system framework that is now popular and well-known:

- System I: fast and intuitive, but prone to cognitive biases and heuristic shortcuts that often lead to demonstrably poor decisions
- System II: slow and deliberate, employing careful analysis to overcome these biases and make rational conclusions

Gigerenzer (1996):²⁰ building on Simon’s bounded rationality, Gigerenzer argues that fast-and-frugal heuristics—simple rules of thumb—are often ecologically rational and has demonstrated that heuristics can outperform complex analytical methods in environments characterized by high uncertainty, sparse data, and time pressure [50, 51].

The Risk Awareness framework [2] emphasized the distinction between Kahneman and Gigerenzer’s work as particularly relevant in flight test. Risk Awareness argued that the domain of uncertainty should determine the most appropriate cognitive approach. System II analytical thinking proves valuable in the risk domain where probabilities are well characterized, while experience-informed heuristics and ‘gut feelings’ often provide better guidance in the pure uncertainty domain where many critical flight test decisions are made.

Klein (1999):²¹ Klein built his Recognition-Primed Decision (RPD) Model based on observations of experienced professionals, *e.g.*, firefighters and military officers, who make rapid and effective decisions without, from his view, explicitly comparing multiple options [52]. RPD involves a two-stage process: an intuitive recognition of a *plausible* course of action based on past experiences (pattern matching) and a mental simulation of that action to ensure its *viability*. Both *plausibility* and *viability* require an implicit assessment of the likelihood of success, a probability assessment. Klein’s RPD model and his description of naturalistic decision-making influenced the Army Field

¹⁸ Daniel Kahneman (1934–2024): Nobel-prize winning Israeli-American psychologist; collaborated with Tversky in the development of prospect theory and established the field of behavioral economics.

¹⁹ Amos Tversky (1937–1996): Israeli mathematical psychologist; collaborated with Kahneman in the development of prospect theory and the field of behavioral economics.

²⁰ Gerd Gigerenzer (1947–): German psychologist.

²¹ Gary Klein (1944–): American research psychologist who worked as a U.S. Air Force psychologist in the mid-1970s.

Manuals on Mission Command and the Army’s formal curriculum on intuitive decision-making [53]. RPD is the motivation for most ‘apprenticeship’ pedagogy, including our Test Pilot Schools, which are fundamentally apprenticeships for the test process (including test point execution, planning, reporting, team building, *etc.*) [54].

The intent of this survey of 20th-century research into decision-making is to underscore that, regardless of model—rational utility, bounded rationality models, slow thinking, fast-and-frugal decisions—decision-making under uncertainty requires an assessment of the expected probability of success (or failure) of an outcome. The better and more self-consistent our probability assessments are, the better risk-informed decisions we will make. We need a systematic method for weighing potential outcomes against their likelihood. We should also employ a common lexicon to harness the wisdom of the crowd to improve the accuracy of our collective judgments. This is the motivation for encouraging the adoption of Bayes Factors in flight test.

III. Bayes Factors

Our knowledge of the world evolves with new evidence and data over time. We are, almost by nature, ‘born Bayesians.’²² The Bayes Factor approach recommended here provides a mathematical framework for combining objective and subjective measures of expectation and for updating our beliefs in light of new evidence and data. This makes it particularly well-suited to the iterative nature of flight test where our understanding evolves continuously from engineering design reviews through bench tests and build-up test campaigns.

The Bayes Factor (BF) was introduced by Jeffreys²³ in his 1939 book *Theory of Probability* [31]. With two complementary hypotheses, the Bayes Factor is the posterior odds of one hypothesis when the prior probabilities of the two hypotheses are equal. More simply, the Bayes Factor, or BF , is the likelihood ratio assessing the relative weight of evidence, data, or judgment between two competing hypotheses. Bayes Factors began to see more widespread adoption by statisticians, scientists, courts, and medical researchers following Kass and Raftery’s 1995 paper on the subject [57].

A direct application of Bayes’ Theorem (Equation A-3) to data is often difficult, because it requires an assessment of the probability of observing a particular piece of data, D (see the example in Section IV.A). Fortunately, the use of Bayes Factors (BF) avoids this step (Appendix A). The Bayes Factor for a particular piece of data or evidence is the likelihood ratio of observing that piece of data under two complementary hypotheses, H and its complement H^c :

$$BF \equiv \frac{P(D|H)}{P(D|H^c)} \quad (2)$$

That is, the Bayes Factor is the ratio of the probability of observing the data or evidence, D , under the assumption that the hypothesis, H , is true to the probability of observing D under the assumption that H is false. Helpfully, any underlying error or unacknowledged assumption affecting

²² Bayesian Brain Theory, a cognitive theory of the mind, proposes that the brain forms Bayesian predictive models of sensory inputs to enable adaptive behavior [55, 56].

²³ Sir Harold Jeffreys (1891–1989), a Fellow of the Royal Society, British geophysicist and statistician.

the data or the accuracy of our judgment is common to both the numerator and denominator in Equation 2, canceling out. This is one of the advantages of using Bayes Factors instead of the conditional probability.

In our flight test risk assessment context, we will use the hypothesis a planned test is safe and the complementary hypothesis that a planned test is unsafe.²⁴ Thus, the Bayes Factor for our flight examples will be

$$BF \equiv \frac{P(D|H_{safe})}{P(D|H_{unsafe})} \quad (3)$$

To use Bayes Factors we have to express our probability in terms of odds ratios, Θ . Although common in gambling, expressing uncertainty in terms of odds is not as common as doing so in terms of probabilities. The odds that a die roll comes up ‘6’ are one-to-five, 1:5, or 5:1 against. That is, a “one-in-six chance” (a probability statement) corresponds to a “one-to-five odds” (an odds statement). It is straightforward to convert probabilities, p , into odds ratios, Θ , by

$$\Theta = \frac{p}{1-p} \quad (4)$$

Odds are converted back into probabilities by

$$p = \frac{\Theta}{\Theta + 1} \quad (5)$$

With the prior odds, $\Theta^{(0)}$, established, the posterior odds, $\Theta^{(1)}$, informed by the data, are

$$\Theta^{(1)} = BF \cdot \Theta^{(0)} \quad (6)$$

Multiple pieces of data, D_i , can be considered simultaneously. The final result is the product of all the individual Bayes Factors for each piece of data, evidence, professional engineering judgment, etc. Thus, the final posterior odds are

$$\Theta^{(1)} = BF_1 \cdot BF_2 \dots \cdot BF_n \cdot \Theta^{(0)} = \prod_i BF_i \cdot \Theta^{(0)} \quad (7)$$

New posterior odds using new Bayes Factors from new data can also be updated over time. This is helpful as additional build-up tests or other data about the system under test become available. The previous posterior odds simply become the new prior odds. Thus

$$\Theta^{(2)} = BF_{n+1} \cdot \Theta^{(1)} \quad (8)$$

Equation 5 is still used to convert the posterior odds back to a probability.

²⁴ Of course, the binary choice of either *safe* or *unsafe* is a false dichotomy, but the use of Bayes Factors requires complementary hypotheses; in practice, we can distinguish a *safe* and *unsafe* test as those that either above or below a threshold for an acceptable level of risk.

IV. Examples

We present four non-test and four flight test examples of using Bayes Factors to make judgments and conclusions. Each of the examples below demonstrates the use of Bayes Factors to aid decision-making under uncertainty, and each example was chosen to highlight different aspects of the application. The examples illustrate the judgment involved in assigning values for Bayes Factors as well as prior odds. We conclude each example with a summary of the lessons learned from the practical application and use of Bayes Factors.

The selection of prior odds and the determination of Bayes Factors will necessarily reflect the subjective nature of probability. This will remain the most philosophically challenging aspect of the framework, and is an inherent and inescapable aspect of decision-making under uncertainty. Our initial beliefs (*priors*) and our assessment of the likelihood ratio of the available evidence (*Bayes Factors*) are not altogether arbitrary. Although subjective, they are informed by our experience, knowledge, and professional judgment. In some circumstances, when the problem is well-defined and the probabilities can be objectively calculated, selecting appropriate priors and Bayes Factors is straightforward. In many real-world applications where uncertainty is present, the choice will be subjective but informed. The flight test examples below include priors that range from 10^{-4} to 20:1. The explicit acknowledgment of our *priors* and our *Bayes Factors* is the basis for meaningful, constructive conversations about our risk assessments. The advantage of quantifying and communicating our judgments and having our assessments open for comparison and scrutiny should be obvious. Additionally, the mechanism for sequentially updating our assessments is built into the Bayes Factor approach as several of the selected examples will also demonstrate.

A. Kahneman and Tversky’s Taxicab Problem

The “Taxicab Problem” is a classic probability puzzle that illustrates the base rate fallacy.²⁵ Amos Tversky and Daniel Kahneman first offered this problem to study subjects in 1972 [58]. They concluded that people underestimate or neglect base rate data leading to incorrect probability judgments. The “Taxicab Fallacy” is also frequently used as an example of a common cognitive shortcut known as the representativeness heuristic.

The Taxicab Problem is this: a hit-and-run accident occurs at night. There are two cab companies in the city operating in distinct green and blue liveries. 85% of cabs in the city are green cabs; 15% are blue cabs. An eyewitness identifies the cab in the hit-and-run accident as blue. The witness is tested under similar conditions and correctly identifies a cab’s color 80% of the time, regardless of the cab’s color. What is the probability that the cab involved in the accident was blue?

Kahneman and Tversky report that for several hundred subjects given slight variations of this question, both the modal and the median response to the question—“what is the probability that the cab was blue?”—was 80%. This is almost double the correct answer of 41%. The common error is that most people focus on the eyewitness accuracy of 80%, concluding that the cab was

²⁵ The cognitive error in which people neglect general statistical information (the base rate) in favor of specific, but less relevant, case-specific information.

blue with 80% probability. This neglects the significant factor that only 15% of cabs in the city are blue—the ‘base rate.’

The correct answer, solved using Bayes’ Theorem (Equation A-3), is

$$P(B|D) = \frac{P(D|B)P(B)}{P(B)P(D|B) + P(G)P(D|G)} = \frac{(0.8)(0.15)}{(0.15)(0.8) + (0.85)(0.2)} = 41\% \quad (9)$$

where B is the hypothesis that the cab was blue, G is the hypothesis that the cab was green, and D is the data—the eyewitness report. Given the complexity of applying Bayes’ Theorem in this form, it is no wonder that most people in the experiment fail to calculate the correct answer.

Solving the problem with Bayes Factors is much simpler than using Bayes’ theorem, Equation 9. The prior odds of a blue taxicab, $\Theta_B^{(0)}$, are 15:85 or 15/85. The Bayes Factor is given directly in the problem by the tested eyewitness reliability of 80%. That is,

$$BF = \frac{P(D|B)}{P(D|G)} = \frac{0.80}{0.20} = 4 \quad (10)$$

and thus the posterior odds of a Blue cab are

$$\Theta_B^{(1)} = BF \cdot \Theta_B^{(0)} = (4) \left(\frac{15}{85} \right) = \frac{12}{17} = 12 : 17 \quad (11)$$

Using Equation 5, we convert these 12:17 posterior odds back into a probability

$$p_B = \frac{\Theta_B^{(1)}}{1 + \Theta_B^{(1)}} = \frac{12/17}{1 + 12/17} = \frac{12}{12 + 17} = 41\% \quad (12)$$

Equation 11 is obviously more straightforward to apply than Equation 9 in solving the Taxicab Problem. This is not unusual. The direct use of Bayes’ Theorem is often challenging due to the need to assess the marginal probability of observing a piece of data. In practice, this requires a self-consistent assessment of the data, the hypothesis, and the complementary hypothesis. Because the Bayes Factor is expressed as a ratio, it sidesteps the need to assess the marginal probability of the data. Furthermore, it is not sensitive to underlying assumptions or biases about the data, since these will cancel out when forming the likelihood ratio. This example offers our first lesson learned on the use of Bayes Factors.

Lesson Learned 1: *BFs are often simpler to apply than direct application of Bayes’ Theorem, producing the “right” answer without needing to assess the problematic marginal probability of observed data (BFs are insensitive to assumptions about the data or evidence and guard against cognitive pitfalls associated with conditional probabilities)*

B. How ‘good’ is that Amazon Review?

Consider an Amazon product review—or hiring reference, restaurant or vacation recommendation, *etc.*—with only five reviews. Four are positive (*e.g.*, ‘five-star’) reviews/recommendations and one is neutral or negative (*e.g.*, ‘two-star’). Given the data, 4-out-of-5 (80%) satisfied customers, what is your expectation of a favorable experience, *i.e.*, a ‘five-star’ experience?

Our Bayes Factor will be the ratio of the probability of observing this data under two hypotheses: that you will have a favorable or an unfavorable experience, $H_{favorable}$ or $H_{unfavorable}$. Under the hypothesis that you have a favorable experience, the probability of having observed the review data is 5/6, because your favorable experience is now a ‘new’ sixth review, five of which are positive (including yours). Likewise, under the hypothesis that you have an unfavorable experience, the probability of observing that data is 2/6. Thus, the Bayes Factor is:

$$BF = \frac{P(D|H_{favorable})}{P(D|H_{unfavorable})} = \frac{5/6}{2/6} = \frac{5}{2} \quad (13)$$

Our posterior odds depend on our prior odds. Let us consider three different individuals: a skeptical consumer who tends to be unsatisfied, a neutral consumer who tends to be equally satisfied and unsatisfied, and an optimistic consumer who tends to have favorable experiences. Each of these consumers will have the same Bayes Factor but different prior odds. For the neutral consumer, there is no *a priori* reason to expect either a favorable or unfavorable outcome, the prior odds are 1:1. This is known as the ‘*uninformed prior*’ (in this case, a uniform distribution). The posterior odds of a favorable experience, given the 80% satisfied customers and a uniform prior, are then

$$\Theta^{(1)} = BF \cdot \Theta^{(0)} = \left(\frac{5}{2}\right) \left(\frac{1}{1}\right) = \frac{5}{2} \quad (14)$$

which corresponds to a 71% probability of a favorable experience for our neutral consumer.²⁶

Consider the optimistic consumer. Let us assume that they generally have twice as many favorable experiences as experiences they regret. The prior odds for this optimistic consumer are 2:1 and the posterior odds are 5:1 ($(5/2)(2/1) = 5/1$). This yields a 83% expectation of a favorable experience. Alternatively, the critical skeptic, who has twice as many unfavorable as favorable experiences, would assign prior odds of 1:2, yielding a posterior odds ratio of 5:4 ($(5/2)(1/2) = 5/4$) and only a 56% expectation of a favorable experience given the 80% positive reviews.

Table 1: Bayes Factors - Amazon Review

	$\Theta^{(0)}$	BF	$\Theta^{(1)}$	$p_{favorable}$	95% CI	$P(p > 0.8)$
Naive Prior	1:1	5/2	5:2	71%	(0.245, 0.843)	35%
Optimistic Prior	2:1	5/2	5:1	83%	(0.587, 0.977)	68%
Skeptical Prior	1:2	5/2	5:4	56%	(0.245, 0.843)	6%

Table 1 summarizes the impact of the various priors. The, perhaps surprisingly, low posterior expectations are the consequence of the relatively sparse data. More data would result in a

²⁶ Reminder, we use Equation 5 to convert the 5:2 odds to a probability: $p = (5)/(5 + 2) = 71\%$.

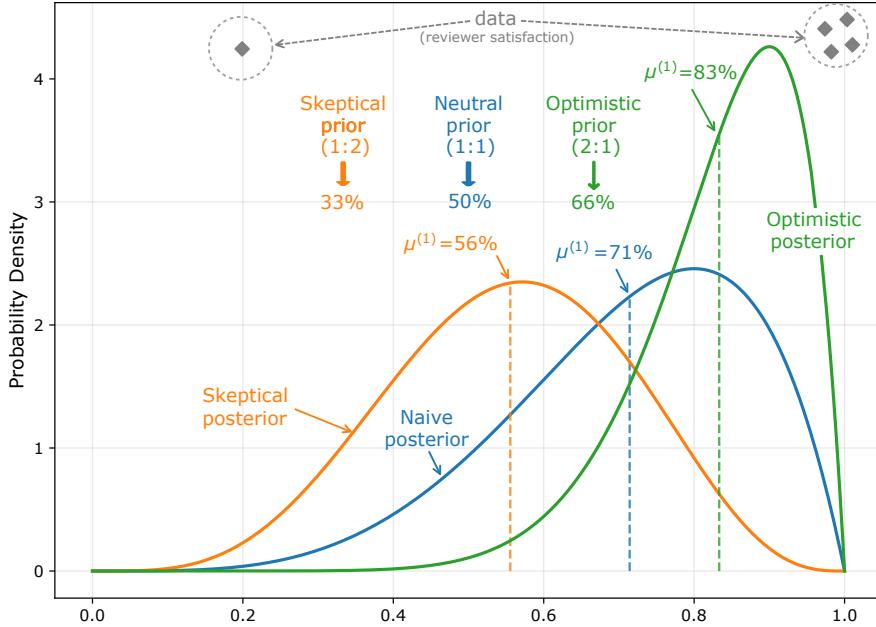


Figure 5: Posterior Expected Satisfaction Given Prior Review Data

more significant Bayes Factor and more confidence in our posterior estimate. Although our posteriors will always depend on our priors, Appendix B demonstrates that increasing the amount of data will ultimately overwhelm the prior odds and the posterior distribution will converge to the ‘true’ distribution.

The posterior distributions in Figure 5 for our three consumers are developed using the methods described in Appendix A. Given a probability distribution, we can establish a 95% confidence interval on our expectation of a favorable experience. We can also answer the question: how likely are our three consumers to have a favorable experience (*i.e.* what is the area under the posterior probability density function for which $p > 80\%$ —80% being taken as an arbitrary threshold for ‘favorable’²⁷). These results are given in Table 1. In general, we will not try to define the posterior probability distribution in our flight test examples, but it is straightforward enough to do so with this example that we have included the result.

This example illustrates that our judgments and conclusions, particularly with sparse data, are heavily influenced by our ‘*priors*.’ With Amazon or restaurant reviews, most consumers would likely guard their expectations with limited data. The same is true with recommendations for hiring decisions, *etc.* The surprisingly low expectation of a favorable experience, stemming from the sparse data in this example, serves as a cautionary example for avoiding the temptation for optimism bias and other planning fallacies. Fat-tailed priors, always useful to account for the possibility of rare, high-impact events, may be even more important in flight-test risk management.

The use of *BFs* also requires us to consider alternative hypotheses explicitly and is thus a good guard against confirmation bias.²⁸ The use of Bayes Factors will not eliminate confirmation

²⁷ This will be analogous to the safe/unsafe dichotomy we will use in our flight test examples; ‘safe’ is intended to mean that an ‘acceptable level of risk’ has been achieved.

²⁸ Confirmation bias—arguably one of the most prevalent and persistent cognitive biases in human affairs—is



Figure 6: Common Final Exam at the United States Air Force Academy

bias, but acknowledging alternative hypotheses is a valuable antidote for avoiding surprises.

Lesson Learned 2: Our judgments and conclusions are strongly influenced by our prior expectations

Lesson Learned 3: The *BF*-framework requires us to acknowledge alternate hypotheses explicitly (and may guard against confirmation bias, planning fallacy, blind spots, and other cognitive errors)

C. Is Charlie a Cheater?

This example is drawn from the author's experience while serving as Chair of the Engineering Division at the U.S. Air Force Academy. A cadet suspected that two of her classmates had cheated on a final exam, having heard them comparing answers to questions during the exam.²⁹ We will call our three cadets *Alice*, *Bob*, and *Charlie*. Alice was sitting in front of Bob and Charlie during a final exam. She reported to the instructor that Bob and Charlie had been talking during the exam and that it sounded like they were comparing answers. Bob and Charlie did admit to talking, but denied cheating. They claimed to have been trying to coordinate a ride to the airport.

Bob would eventually admit to having cheated off Charlie, but claimed that Charlie had no knowledge of Bob's cheating. Interestingly, Bob received a better grade, missing only eight questions, compared to Charlie's 13 incorrect answers on the 50-question exam. However, of the eight questions they jointly missed, their incorrect answers were all identical. Bob admitted to cheating and started honor probation immediately. Charlie had a previous honor violation. A second honor conviction would likely result in his dismissal, so he clearly had an incentive to deny cheating. Ultimately, Charlie was found not guilty at an Honor Board of his peers. He appealed his academic penalty of zero points on the exam, based on the not-guilty finding. Charlie's appeal

widespread and consistently replicated in studies; it has a direct implication in decision-making under uncertainty: once an individual makes a decision, there is a strong and persistent tendency to look for confirmatory information that supports the decision.

²⁹ The Academy takes its Honor Code quite seriously; lying, stealing, or cheating are grounds for dismissal.

came to the author as the Chair of the Engineering Division for adjudication and final decision. So the judgment under uncertainty the author had to make: is Charlie a cheater?

The relevant data in the case includes both exonerating and incriminating evidence. This includes:

1. Alice observed Bob and Charlie talking
2. Bob and Charlie admitted to talking, but claimed to be arranging a ride to the airport
3. Charlie was found “Not Guilty” by an Honor Board of his peers
4. Charlie’s prior honor conviction provides a clear incentive to lie about cheating
5. Bob admitted to cheating off Charlie, so there is a *prima facie* case that cheating did occur
6. No other cadets sitting around Alice, Bob, or Charlie reported anything suspicious; under the Academy’s Honor Code, cadets may not tolerate cheating and are obligated to confront suspicions of violations (the other cadets sitting around Bob and Charlie were also friends of Bob and Charlie)
7. There were two versions of the exams, A-versions and B-versions, alternating seat-by-seat. Bob and Charlie, sitting in adjacent seats, should have had different versions, but someone had switched the order before the exam started so that both Bob and Charlie had the same version of the exam
8. The similarity and correlation of incorrect answers between Bob and Charlie’s exams can be compared with all 186 cadets taking the same exam version

We start with two complementary hypotheses: H_0 , Charlie is innocent; H_1 , Charlie cheated. The Bayes Factors will be the relative likelihood of observing the data given above under these two hypotheses:

$$BF_i = \frac{P(D_i|H_1)}{P(D_i|H_0)} \quad (15)$$

Table 2 summarizes the evidence as well as the BF s the author used in judging Charlie’s likelihood of having cheated. BF_1 in case of Alice’s report, assuming 80% eyewitness reliability, is 4:1, just as it was in the prior taxicab example (Section IV.A). The author gave Bob and Charlie the benefit of the doubt regarding their claim that they had only been arranging a ride to the airport. Assessing equal probability of observing the talking under either hypothesis, $BF_2 = 1:1$.

The “not guilty” verdict from the Honor Board is not as exonerating as it might otherwise be, given a persistent low confidence in the honor system among cadets at the time. The low confidence among cadets and the open admission by some cadets that they would never convict a fellow cadet make this line of evidence less compelling. Nevertheless, this was the basis for Charlie’s appeal of his academic penalty, and it deserved careful consideration. The author ultimately assessed the probability that Charlie was innocent given the “not guilty” verdict as 75%, yielding a $BF_3 = 1:3$.

Table 2: Is Charlie a Cheater?

Evidence	BF_i	Comment
1 observed by student talking	4/1	$\approx 80\%$ eyewitness reliability
2 admitted to talking	1/1	gives Bob & Charlie the benefit of the doubt
3 “not guilty” at honor board	1/3	cadet & faculty confidence in the honor system
4 prior honor violation	4/3	slight incentive to lie; some benefit of the doubt
5 Bob’s admission of cheating	3/2	acknowledgment that cheating did occur
6 no reports from other cadets	4/5	surrounded mostly by friends
7 duplicate adjacent exams	5/3	exam versions were deliberately exchanged
		$P(x_8 H_1) = 1$
8 exam similarity	958/3	$P(x_8 H_0) = 3/958$ (only 3/958 students with perfect incorrect-answer correlation)

Charlie’s previous honor violation and clear incentive to deny cheating is incriminating, but the author again gave Charlie a slight benefit of the doubt and assigned a relatively weak $BF_4 = 4:3$ to this data. Bob’s admission to cheating off Charlie, allegedly without Charlie’s knowledge, was dubious (particularly given the fact that Bob performed significantly better on the exam than Charlie). Still, the author also assigned a relatively weak $BF_5 = 3:2$ to this data. The absence of reports from other cadets in the area, all of whom were friends with Bob and Charlie, was not compelling and received a weak, exonerating $BF_6 = 4:5$. The deliberate exchange of alternate versions of the exam was incriminating, though Bob and Charlie denied any part in the exchange. The author judged this to be mildly incriminating and assigned a $BF_7 = 5:3$.

All of the previously considered BF s in this example have been qualitative judgments given in quantitative ratios. In the case of the data on the similarity between Bob and Charlie’s exams, we can be precise and quantitative. A total of 186 students took the final exam. Figure 7 depicts a heatmap of the 186x186 pairwise dot products of incorrect answers of every copy of the same version of the exam. This measure of the correlation coefficient between exams shows how similar students’ incorrect answers were. A histogram of this data shows that 958 pairs of exams were correlated with eight or more incorrect answers (there were 50 questions on the exam). In only three cases—of 958 pairs of exams with eight or more incorrect answers—did every single incorrect answer correlate exactly with the other exam. Bob and Charlie were one of these three pairs of perfectly correlated incorrect answers. The probability of observing this perfect correlation of incorrect answers, under the assumption that Charlie is a cheater, is 1. The probability of observing this data under the assumption that Charlie is innocent is 3/958 given the class-wide data. Therefore, the author assessed a Bayes Factor of 958/3 for this piece of evidence.

The combined Bayes Factor is then

$$BF = \left(\frac{4}{1}\right) \left(\frac{1}{1}\right) \left(\frac{1}{3}\right) \left(\frac{4}{3}\right) \left(\frac{3}{2}\right) \left(\frac{4}{5}\right) \left(\frac{5}{3}\right) \left(\frac{958}{3}\right) = 1135 \quad (16)$$

Estimating the prevalence of cadet cheating at 1% gives prior odds of a cadet cheating of 1:99. The

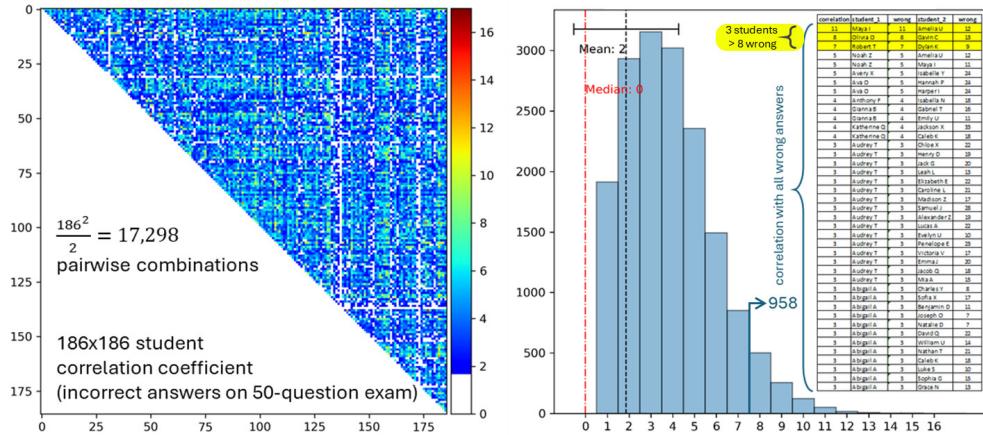


Figure 7: (L) 186x186 pairwise correlation coefficient between all 186 students who took the final exam; (R) histogram of the 17,298 pairwise combinations between all student exams

posterior odds and probability that Charlie cheated in the final exam are then

$$\Theta^{(1)} = BF \cdot \Theta^{(0)} = \left(\frac{1135}{1} \right) \left(\frac{1}{99} \right) = 11.5 \rightarrow 92\% \text{ Charlie cheated} \quad (17)$$

On the basis of this Bayesian reasoning, the author denied Charlie's request and sustained the academic penalty he had received on the Final Exam.

This example illustrates two advantages of using Bayes Factors. First, the assigned *BFs* indicate how strong a particular piece of evidence is judged to be. Without the exam similarity, we would generously assess the probability that Charlie cheated as less than 5%. However, the extremely strong *BF* of more than 300 indicates how incriminating this piece of evidence is compared to the “not guilty” verdict from the Honor Board.³⁰ Secondly, this example shows the opportunity for constructively communicating judgment. The author was able to acknowledge and incorporate all the exonerating evidence Charlie presented. The author was also able to explain his reasoning to Charlie, along with the relative weight of the various pieces of evidence. Charlie accepted the rationality of the author’s conclusions and judgments with dispassionate equanimity. As we will see in the examples of safety and technical risk discussions below, Bayes Factors provide a common framework for a meaningful discussion of the relative strengths of conflicting data, and give us a systematic way to communicate our judgments and conclusions.

Lesson Learned 4: *BFs* can mix both subjective and objective judgments

Lesson Learned 5: The *BF*-framework decomposes complex data sets into separate chunks and reveals the relative strength of different pieces of evidence (combining multiple pieces of evidence is as simple as multiplying *BFs*)

Lesson Learned 6: *BFs* provide a systematic structure for meaningful conversations and communicating our experience, judgment, and decisions

³⁰ Relying on the exam similarity alone and concluding that the probability Bob and Charlie had not cheated was 3/958 would be erroneous; an example of the prosecutor’s fallacy (see Section V). With the same prior odds of 1:99, the evidence from the exam similarity alone yields only 3:1 posterior odds that Charlie cheated.

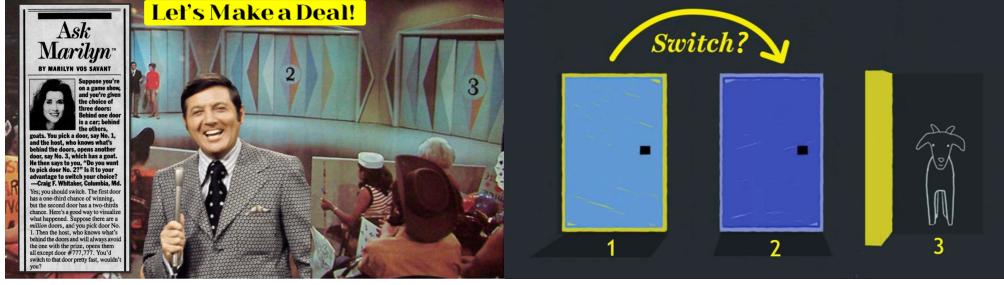


Figure 8: The Monty Hall Problem. Do you switch doors?

D. Monty Hall Problem

This is another classic probability question [59, 60]. Although not new at the time, it went viral in 1990 when Marilyn vos Savant shared it in her Parade Magazine column [61]. The column received thousands of letters from readers who disagreed with her answer. Many of the letters were from “experts” with PhDs; some of them were quite sexist.³¹

The Monty Hall problem is named for the host of the game show, *Let’s Make a Deal* (Figure 8), in which the contestant is asked to choose between one of three doors: Door 1, Door 2, or Door 3. Behind one of the doors is a car. There are goats behind the other two. Suppose a contestant selects Door 1. The host, who knows what is behind each door, opens Door 3, revealing a goat. The host then gives the contestant the opportunity to change their choice of doors to Door 2. The somewhat counterintuitive solution, and the one that raised the ire of all the letter writers, is that switching doors doubles the contestant’s chance of winning the car. The contestant, who originally had a $1/3$ chance of being right with Door 1, has a $2/3$ chance of winning by switching to Door 2.

Let H_i be the hypothesis that the car is behind Door i . We will refer to the data of the goat behind Door 3 as G_3 . As with the Taxicab Problem (Example IV.A), proving that switching doors doubles the contestant’s chances of winning via Bayes’ theorem is mildly torturous:

$$\begin{aligned} P(H_2|G_3) &= \frac{P(G_3|D_2)P(D_2)}{P(D_1)P(G_3|D_1) + P(D_2)P(G_3|D_2) + P(D_3)P(G_3|D_3)} \\ &= \frac{(1)(1/3)}{(1/3)(1/2) + (1/3)(1/1) + (1/3)(0)} = \frac{1/3}{1/6 + 2/6} = \frac{2}{3} \end{aligned} \quad (18)$$

This gives the correct answer but no insight. The Bayes Factor method is more direct and intuitive. Consider the odds comparing Door 2 to Door 1. The prior probability for winning with either door is $1/3$, so the prior odds ratio for Door 2 to Door 1 is 1:1. The Bayes Factor comparing the likelihood of winning with Door 2 to Door 1 given the goat behind Door 3 is

$$BF_{D_2:D_1} = \frac{P(G_3|H_2)}{P(G_3|H_1)} = \frac{1}{1/2} = \frac{2}{1} \quad (19)$$

The numerator—the probability of being shown the goat behind Door 3 given that the car is behind

³¹ She published many of the letters in a subsequent column, with only a trace of gloating [62].

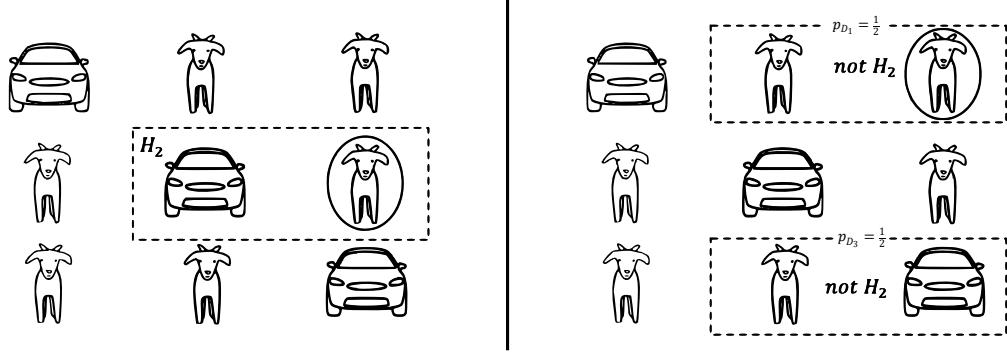


Figure 9: (L) H_2 : hypothesis that the car is behind Door 2;
(R) $\text{not } H_2$: hypothesis that the car is not behind Door 2

Door 2—is 1 because Monty can only show you the goat behind Door 3 (Figure 9). The probability in the denominator is 1/2 because if the car is behind Door 1, Monty can show you a goat behind either Door 2 or Door 3, so the probability of seeing G_3 is 1-out-of-2. With the Bayes Factor from Equation 19, the posterior odds for winning with Door 2 compared to Door 1 are

$$\Theta_{D_2:D_1}^{(1)} = BF_{D_2:D_1} \cdot \Theta_{D_2:D_1}^{(0)} = \left(\frac{2}{1}\right) \left(\frac{1}{1}\right) = \frac{2}{1} \rightarrow 66\% \text{ winning by switching to Door 2} \quad (20)$$

The Bayes Factor approach is a far simpler solution to the Monty Hall problem.³²

To make things more interesting, and reflecting the definition that *probability* is our subjective relationship with uncertainty, suppose the contestant has additional data, considerations, or suspicions about Monty’s behavior. *E.g.*, consider the case that, through repeated observation of *Let’s Make a Deal*, the contestant concludes the car is more often behind Door 1. Or perhaps the contestant suspects that the producers do a poor job of randomizing doors and never repeat the same door week to week. Or perhaps we suspect that Monty knows about the solution to the Monty Hall problem, will only offer the contestant the chance to switch if they initially selected the correct door. These are all subjective judgments or observation-informed expectations about the real world that the contestant can systematically quantify using Bayes Factors.

³² The solution by Equation 20 was simplified by recasting it in terms of the odds ratio comparing Door 2 to Door 1. A brute force approach that directly solves for the posterior odds of Door 2 given G_3 is also possible. The prior odds of a car behind Door i are $\Theta_i^{(0)}$. Thus, $\Theta_2^{(0)} = 1 : 2$. The Bayes Factor for switching doors is

$$BF = \frac{P(G_3|H_2)}{P(G_3|H_2^c)} = \frac{P(G_3|H_2)}{\frac{1}{2}P(G_3|H_1) + \frac{1}{2}P(G_3|H_3)} = \frac{1}{\left(\frac{1}{2}\right)\left(\frac{1}{2}\right) + \left(\frac{1}{2}\right)(0)} = \frac{1}{1/4} = 4$$

The probability of observing G_3 under the assumption of H_2 being true is 1, as before. The probability of observing G_3 under the assumption of H_2^c must consider both H_1 and H_3 (Figure 9). This gives 1/4 for $P(G_3|H_2^c)$. The posterior odds for switching doors are then

$$\Theta_2^{(1)} = BF \cdot \Theta_2^{(0)} = \left(\frac{4}{1}\right) \left(\frac{1}{2}\right) = 2 \rightarrow 66\% \text{ winning by switching to Door 2 as before}$$

The brute force approach is not as simple and elegant as Equation 20, but it is still simpler and more intuitive than the conditional probability calculation required for the Bayes theorem approach, Equation 18.

H_0 : unsafe test	Prior Estimate Safe	Bayes Factor (Likelihood Ratio):
H_1 : safe test	Test : 100:1 against	$\frac{P(\text{Evidence} \text{SAFE Test})}{P(\text{Evidence} \text{UNSAFE Test})}$
	Prior OR = $\frac{1}{100}$	
Evidence	$\frac{BF}{\Delta}$	
① Certified Store (F-16)	① 2:1	$\frac{2 \cdot 3 \cdot 2 \cdot 1 \cdot 1 \cdot 2 \cdot 1 \cdot 2 \cdot 1}{1 \cdot 1 \cdot \frac{1}{3} \cdot 1 \cdot 1 \cdot 1 \cdot 1 \cdot 3} = 32$
② Similar Store	② 3:1	post OR = $(32) \frac{1}{100} \approx \frac{1}{3}$
③ CFD Flow Field	③ -	$P(\text{SAFE Test}) = \frac{\frac{1}{2}}{1 + \frac{1}{3}} = \frac{1}{2} = 50\%$
④ Wind Tunnel Safe Separation	④ -	⑦ AFSES Judgment 3:1 \Rightarrow post OR: 1:1
⑤ Drop Test (Pyram Lads)	⑤ -	ISG maturity (IS) $(3)(\frac{1}{3}) = 1$
⑥ Post Sep Dynamics < 6.025	⑥ -	New trajectory further west?
⑦ Integration Test	⑦ 1:1	Yanks engineers confidence
⑧ Detailed System Understanding	⑧ 1:1	⑧ Post-sep control authority 2:1
⑨ Measured Mass Properties	⑨ 2:1	Posterior OR: 2:1
⑩ Autopilot Control Authority	⑩ -	$\Rightarrow P(\text{SAFE Test}) = 66\%$
⑪ Trajectory Monte Carlo simulation	⑪ 10:1	
⑫ Separation Test Article	⑫ -	
⑬ Flight Termination System	⑬ 1:5	
⑭ Trackers	⑭ -	
⑮ Stability Analysis	⑮ 2:1	
⑯ Difference S & 20°, reduced polar	⑯ 1:3	

Figure 10: Project Serene TAB/SRB notes showing “back of the envelope” Bayes Factor calculations

Lesson Learned 7: BFs can account for and incorporate other factors that were not part of the original problem statement

E. Project Serene

Project Serene was a non-traditional, long-range weapon test with considerable unknowns and uncertainty. The weapon had a limited pedigree of previous employment, but the platform and method of employment under test were entirely new. For various programmatic reasons, many of the typical risk reduction and build-up tests were not possible for the project. These were still listed as potential evidence (factors), capturing how the posterior odds could be influenced by additional data, but were not assigned *BFs* in computing the posterior odds. The initial data and initial planning *BFs* are given in Table 3.

The detailed rationale for each *BF* in Table 3 is not as important as the technique. Figure 10 contains the author’s notes from the initial planning meeting. This demonstrates how straightforward it is to track multiple lines of evidence during a meeting. Following extended discussion regarding each factor with all the experts at the table, the author recorded his real-time assessment of the Bayes Factor, which was the likelihood of that factor under the hypothesis that the test would be safe relative to that factor under the hypothesis that the test would be unsafe. Overall, the approach was a straightforward *piece-of-scratch-paper* technique that reveals the author’s judgment of the relative strength of the different factors under consideration.

The overall initial Bayes Factor from the initial planning meeting for the factors in Table 3 was 32. Given the weapon had some record of previous employment, the author granted a “generous” prior odds ratio for a safe test of 1:100. With this prior, the posterior odds were

Table 3: Project Serene - Initial Bayes Factors

Evidence	BF	Comment
1 certified store	4:1	but different platform
2 similar store	3:1	engineering data supports similarities
3 CFD flow field predictions	—	not accomplished
4 wind tunnel safe separation test	—	not accomplished
5 drop test (pylon loads)	—	not accomplished
6 post-sep dynamics (6-DOF analysis)	—	not accomplished
7 integration testing	1:1	limited scope
8 detailed system understanding	1:1	lack of engineering data & technical detail
9 measured mass properties	2:1	somewhat limited
10 'autopilot' control authority	—	unknown
11 Monte Carlo trajectory analysis	10:1	acceptable; unfavorable edge cases indicate potential off-range impact
12 separation test article	—	not possible
13 flight termination system	1:5	not feasible
14 trackers	—	not feasible
15 stability analysis	2:1	statically stable article
16 reversed pylon installation	1:3	uncertainty from differences

$$\Theta^{(1)} = BF \cdot \Theta^{(0)} = \left(\frac{32}{1}\right) \left(\frac{1}{100}\right) \approx \frac{1}{3} \rightarrow 25\% \text{ (expectation of safe test)} \quad (21)$$

During the planning meetings, the author also routinely solicited *BFs* from the other meeting participants. A significant range of opinions emerged regarding the safety of the proposed test. The lack of expert consensus provided a direct measure of the overall uncertainty, underscoring a wide confidence interval for possible outcomes.

The author's initial assessment of a posterior odds of 1:3, or 25% probability of a safe test, prompted a request for additional information. In particular, the author, serving as the Test Acceptance Authority for Project Serene, requested expert opinions and judgments from the Air Force SEEK EAGLE Office (AFSEO). AFSEO provided further analysis along with their judgment. This additional data yielded an updated Bayes Factor (Table 4).

Table 4: Project Serene - Bayes Factors II

Evidence	BF	Comment
17 AFSEO judgment	3:1	lack of engineering data results in significant residual uncertainty; store similarity alleviates some risk based on engineering judgment

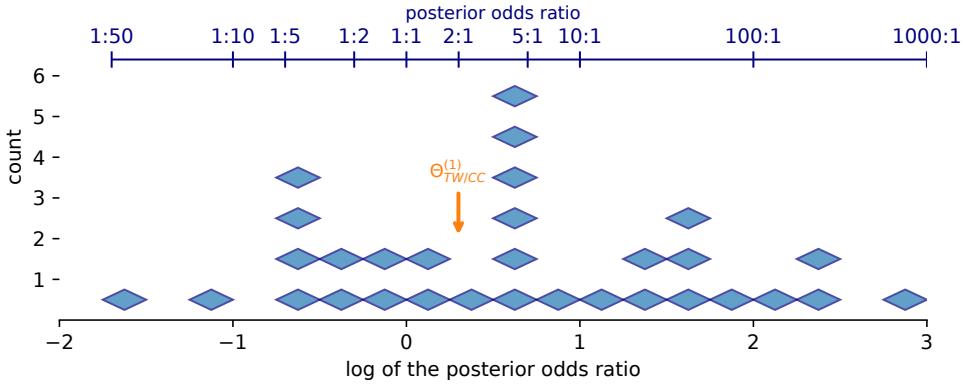


Figure 11: Project Serene posterior odds ratios for all 30 safety review participants

The updated posterior odds given AFSEO’s recommendation are given by the product of the previous posterior odds (the ‘new prior’) and the new Bayes Factor:

$$\Theta^{(2)} = BF_2 \cdot \Theta^{(1)} = \left(\frac{3}{1}\right) \left(\frac{1}{3}\right) \approx 1 \rightarrow 50\% \text{ (expectation of safe test)} \quad (22)$$

Given even posterior odds following AFSEO’s updated analysis, the author was still unwilling to approve the test. The most significant residual concerns were due to uncertainty over the post-separation control authority of the system under test and the kinematic potential for the weapon to depart the range and impact a populated area. This prompted a request for additional engineering data from the vendor and discussions with the engineering team. The additional evidence resulted in a third Bayes Factor (Table 5).

Table 5: Project Serene - Bayes Factors III

Evidence	BF	Comment
18 post-separation control authority	2:1	some confidence in post-separation dynamics & updated kinematic analysis

As before, updated posterior odds considering the new data are calculated by the product of the previous posterior odds (the ‘new’ prior) and the new Bayes Factor:

$$\Theta^{(3)} = BF_3 \cdot \Theta^{(1)} = \left(\frac{2}{1}\right) \left(\frac{1}{3}\right) \approx \frac{2}{1} \rightarrow 66\% \text{ (expectation of safe test)} \quad (23)$$

At this point, after several months of planning, the author asked all participants (approximately thirty aircrew, test safety, range safety, test engineers, project planners, program engineers, and technical advisors) to give their own final, individual posterior odds ratio that the test would be successful. The results, Figure 11, reveal almost five orders of magnitude in differing opinions. The author was ultimately “more-confident-than-not” (66%:33%) that the test would be safe and approved the test (the author’s $\approx 2/3$ confidence is indicated in Figure 11).



Figure 12: B-21 First Flight in the fall of 2023

Differences of opinions in conclusions are either due to different *BFs*, *i.e.*, the relative likelihood of the data under different hypotheses, or due to different prior odds. By framing the question in terms of Bayes Factors and articulating our numbers, we can expose those differences in a meaningful way. This is further explored in Example IV.H.

Ultimately, the first launch of Project Serene was not successful due to a technical deficiency, although the overall test was not unsafe. The second attempt following additional work was successful. A single example is not definitive, but the author finds the consistency between the Bayesian reasoning in this example, the range of test planning expectations, and the eventual outcome to be consistent and representative.

Lesson Learned 8 : *BFs* are straightforward to perform real-time and via “back-of-the-envelope” calculations

Lesson Learned 9: *BFs* provide a common framework to express and expose the full range of subject matter expert opinions; the range of opinions provides a good estimate for uncertainty and for bounding judgments with confidence intervals

F. First Flights: Northrop Grumman B-21 Raider & Hermeus Quarterhorse Mk1

The Bayes Factors technique for articulating stakeholder judgments and communicating risk assessments with a common lexicon proved useful through the ground test build-up and first flight readiness reviews (FFRR) for several programs at Edwards AFB over the past couple of years. This includes approval of the Northrop Grumman B-21 Raider first flight, approval of the Hermeus Quarterhorse Mk1 first flight, and early test planning for the General Atomics YFQ-42A and Anduril YFQ-44A Combat Collaborative Aircraft (CCA) programs. These programs also illustrate the range of *priors* used by the author.

During the First Flight Readiness Review for the B-21 in 2023, the author used Bayes Factors to assess technical, programmatic, and security risks in addition to safety risks. Although we



Figure 13: Hermeus Quarterhorse Mk1 undergoing taxi testing at Edwards AFB

will not provide details for security reasons, the application was straightforward and incorporated data from ground testing, taxi testing, aircrew and test team training, mission rehearsals, airworthiness, and engineering reviews. Sequential updates to the posterior odds were made over the course of several weeks, providing a useful framework for communicating judgments and opinions. The author's initial prior probability for a safe first flight for the B-21, before considering any engineering analysis, airworthiness data, test results, *etc.*, was 10^{-4} . This low prior reflected the magnitude of evidence the author wanted to have to provide confidence in an acceptable level of risk. The Bayes Factors for the engineering evidence, team readiness, *etc.* easily overwhelmed the skeptical prior odds and the B-21 Combined Test Force executed a flawless first flight.

Hermeus, a relatively new aerospace and defense technology company, was founded to develop and build high-Mach and hypersonic aircraft capable of reaching speeds up to Mach 5. The initial effort in a series of planned aircraft is the Quarterhorse program. Hermeus' Quarterhorse Mk1, an uncrewed, remotely piloted aircraft powered by a single GE J85 turbojet engine, had a successful first flight at Edwards AFB in May 2025 following a year of rapid development and fast-paced ground testing.

As with the B-21 first flight, the author relied on the Bayes Factor technique to assess the Mk1 platform and test team readiness for first flight. For the Mk1 program, the author assessed a prior odds of $1 : 10^2$ before consideration of any evidence. The difference in prior odds between the Mk1 and B-21 first flight efforts reflects a difference, based on the author's judgment, in the appropriate threshold for an acceptable level of risk. B-21 is a manned platform with significant program consequences should the first flight be unsafe; Mk1 was an unmanned, rapid prototype that was expected to operate within well-defined and constrained limits in R-2515 restricted airspace at Edwards AFB.

As with other examples, the use of Bayes Factors proved helpful in articulating judgment and discussing the strength of evidence. Bayes Factors also proved useful in making technical, programmatic, and schedule risk assessments and in communicating the author's opinion and judgment.

Lesson Learned 10: *BFs are applicable and useful with other probabilistic risk assessments including technical, programmatic, and schedule risk*

G. Aeolus Aero RPTV-1

Aeolus Aero³³ is a new aerospace company developing a prototype test vehicle under a Cooperative Research and Development Agreement (CRADA) with the Air Force Research Laboratory. Aeolus' initial effort, the Rapid Prototype Test Vehicle 1 (RPTV-1), is an uncrewed, remotely piloted aircraft powered by a single GE J85 turbojet engine. In late spring 2025, the RPTV-1 had a successful first flight at Edwards AFB after a sequence of build-up ground testing.

Table 6: Aeolus RPTV-1 initial planning meeting Bayes Factors

Evidence	BF	Comment
1 initial taxi demos	3:1	surrogate taxi test of Aeolus C2 architecture
2 taxi testing (Edwards)	—	to be accomplished
3 comm-link range testing	—	to be accomplished
4 wind-tunnel stability data	5:1	good results (USAFA's subsonic wind tunnel)
5 lack of autopilot stabilizing	1:4	unknown initial trim settings
6 FTS	1:1	hybrid design, need additional info
7 glide-cone footprint	10:1	kinematic constraint (good risk mitigation)
8 lakebed landing option	5:1	reduces risk exposure
9 datalink transfer to alt mode	1:3	somewhat limited; need better understanding
10 engine runs	—	need x-wind data for limits
11 airworthiness	—	pursuing alternate pathway with FAA

When Aeolus' leadership first approached the 412th Test Wing (TW) to discuss testing the RPTV-1 in R-2515 restricted airspace, they expressed a “tolerance for risk” and desire to “move fast and learn quickly.” The RPTV-1 program was assigned to the 412 TW’s Experimental Test Force (ETF) and given the mandate to support the Aeolus First Flight test campaign. Bayes Factors proved to be a useful method for continuously communicating with Aeolus’ leadership on the author’s opinion of the readiness of the RPTV-1. Table 6 includes the Bayes Factors the author used at the very first, executive kickoff meeting with Aeolus leadership. The author included placeholders for Bayes Factors of testing that would be accomplished when the RPTV-1 arrived at Edwards, *e.g.*, taxi testing and command-and-control link range testing, and would be part of the eventual first flight approval. These served as useful “good-faith” markers and helped establish a common understanding of expectations during initial planning.

Following taxi testing in the fall of 2024, the initial Bayes Factor for the evidence considered at the Test Approval Board (TAB) is given in Table 7. The author assigned a prior odds of 1 : 1000, reflecting the author’s judgment of the acceptable level of risk. Although the author respected and supported the Aeolus Team’s tolerance of elevated risk, and the government had no equity at risk in the RPTV-1, the author could not accept risk to other Edwards assets.³⁴ The *BF* of 20,000 implied by Table 7 combined with a prior odds ratio of 10^{-3} is

³³ A pseudonym

³⁴ Put crudely: it was ok if RPTV-1 crashed on its first flight, but it was not acceptable for it to crash into the B-21 compound.

Table 7: Aeolus RPTV-1 Test Approval Board

Evidence	BF	Comment
1 Team & Engineering	10:1	impressive, small team with trusted pedigree and demonstrated experience & competence
2 successful taxi testing	4:1	good results thru high-speed taxi
3 datalink margins	2:1	demo sufficient link margin & redundancy
4 design & certified parts	1.5:1	qualification of critical design elements
5 lack of stab augmentation	1:1.5	residual design uncertainty; initial trim settings not validated
6 footprint containment	10:1	kinematically contained to remain within the boundaries of the lakebed
7 loss link logic	5:1	reasonable engineering design; redundant architecture
8 FTS	5:1	independent, but non-redundant fuel-cutoff system; components qualified & extensively tested
9 airworthiness	—	pursuing Title 14 exemption

$$\Theta^{(1)} = BF \cdot \Theta^{(0)} = (20,000) \left(\frac{1}{1000} \right) = 20 \rightarrow 95\% \text{ (expectation of safe test)} \quad (24)$$

The author’s risk assessment of 95% expectation of a safe test is based solely on the author’s subjective judgment, capturing both the prior expectations and the author’s opinion of the relative strength of the factors considered. This approach helped quantify and articulate the most significant factors, providing a transparent basis for meaningful, technical discussions between team members and the test approval authority.

The completion of taxi testing and test approval coincided with the arrival of winter precipitation at Edwards AFB and the closure of the lakebeds. One of the risk mitigations for approval of RPTV-1 operations was a restriction to lakebed-contained operations. The closure of the lakebeds resulted in an extended program delay. During the three-month delay between test approval and the ultimate first flight, several key test personnel left the RPTV-1 program, driving a follow-on, FINAL First Flight Readiness Review (FFRR). The factors considered are given in Table 8.

The preponderance of factors considered during the FINAL First Flight Readiness Review indicated an increased risk relative to the earlier assessment. The previous Test Acceptance Board posterior odds of 20:1 formed the basis for the prior odds at the readiness review. The intent of the FINAL readiness review was to consider the impacts to overall risk given the change in circumstances, personnel, currency, *etc.* The *BFFFRR* from Table 8 of 1/6th, the final posterior odds are

Table 8: Aeolus RPTV-1 FINAL First Flight Readiness Review (FFRR)

Evidence	BF	Comment
10 long lay off	1:3	3-mo break due to red lakebeds
11 personnel turnover	1:4	departure of key Aeolus personnel (reduction to previous strong credit)
12 mitigation	2:1	observed/reviewed test conduct to build confidence in team readiness
13 FAA airworthiness waiver	1:1	waiver given; no analysis/judgment implied
14 team currency	2:1	mission rehearsals, training, EP sims
15 company/programmatic pressure	1:2	unavoidable... potential for “Drift”

$$\Theta^{(2)} = BF_{FFRR} \cdot \Theta^{(1)} = \left(\frac{1}{6}\right) \left(\frac{20}{1}\right) = \frac{10}{3} \rightarrow 77\% \text{ (expectation of safe test)} \quad (25)$$

Although this represents an increased risk (in the author’s judgment) to the earlier test approval, the author was sufficiently confident given this analysis to approve the test. The Aeolus team successfully completed the first flight of RPTV-1.

As with Project Serene (Section IV.E), the iterative use of Bayes Factors was useful in continuously updating our judgment based on new information. This is the essence of ‘*Being Bayesian*.’ The systematic record of Bayes Factors from earlier meetings was also helpful in recalling the basis of our previous judgments (*e.g.*, appropriately discounting the earlier credit given for trust in key personnel following their departure from the current test effort) and for quantifying and communicating our judgments.

Lesson Learned 11: are useful for tracking changing judgments over time and for communicating those judgments with stakeholders; includes communicating progress towards an objective and the relative importance of factors for “getting to yes” (avoiding a sequence of “bring me a rock” drills^a)

^a “Bring me a rock” is a popular metaphor describing leaders who have a specific idea in mind but fail to communicate it clearly

H. C-17 Missile Defense Agency High-Altitude Airdrop

The air-launched intermediate-range ballistic missile (AL-IRBM) is a target missile developed by the Missile Defense Agency (MDA) to test the Ballistic Missile Defense Shield (BMDS). The AL-IRBM is airdropped and launched from a C-17A to simulate an Intercontinental Ballistic Missile (ICBM) trajectory. After launch, the AL-IRBM is targeted and intercepted using the missile defense shield. The high-altitude (25k ft) airdrop from the C-17 requires deliberate aircraft depressurization, exposing crew members to physiological risks of hypoxia and decompression sickness (‘the bends’). The danger is compounded by the remote operational areas required for



Figure 14: 418 FLTS C-17 taking off from Edwards AFB & deploying an air-launched medium-range ballistic missile

the mission, with the aircraft typically operating more than five hours away from medical facilities equipped with hyperbaric chambers. Delaying treatment for decompression sickness (DCS) beyond four hours significantly increases the risk of catastrophic adverse outcomes for DCS.

In 2023, a 33-year-old civilian employee of the Missile Defense Agency died of cardiac arrest after experiencing decompression sickness during an Air Mobility Command (AMC) C-17 high-altitude airdrop mission in Alaska. Following this Class A mishap, responsibility for all future high-altitude C-17 MDA missions was returned to Air Force Test Center, and the 418 FLTS was assigned Participating Test Organization responsibility.³⁵

During regular review and approval of the Test and Safety planning package for MDA missions in 2025, the risk of DCS, one of four long-standing test hazards in the safety package, was reassessed as HIGH RISK due to new data from the Air Force's 711 Human Performance Wing (HPW). Given almost two decades of experience with similar airdrop missions, the majority of which had previously been designated MEDIUM RISK, there was substantial disagreement among subject matter experts on the Safety Review Board (SRB) regarding the appropriate level of risk. An approach using Bayes Factors proved to be a useful framework for building a common lexicon and understanding the various perspectives.

The SRB considered two different sets of data during the safety assessment: previous MDA test program flight physiological data and DCS research from the 711 HPW and the dive community. From 24 previous MDA test campaigns since 2004, with an average of 4-5 flights per program and approximately 20 personnel per flight, the SRB estimated there were 2000 previous person-flight hours of high-altitude exposure (approximately one hour per person per mission). During these 2000 previous person-flight hours, there were seven reports of physiological incidents: six involving hypoxia or hyperventilation and the MDA civilian fatality while being treated for DCS.

The USAF's Altitude DCS Risk Assessment Computer (ADRAC) predicted a 1% chance of DCS for the test altitude and exposure time. The SRB's independent reviewer from Aerospace Flight Medicine agreed with ADRAC's estimate of a 1% chance of DCS and added the expectation of a 1% chance of a catastrophic outcome for any given case of DCS based on a comprehensive review of data from the diving community. The probability of a catastrophic event is then estimated

³⁵ A Participating Test Organization (PTO) provides resources and capabilities to support a test conducted by another, lead test organization, in this case, the Missile Defense Agency.

		Mishap Severity			
		Catastrophic ~ 1 (Class A)	Critical ~ 2 (Class B)	Marginal ~ 3 (Class C)	Negligible ~ 4 (Class D/E)
Probability of Occurrence	Level A ~ "Frequent" > 10%	HIGH	HIGH	MEDIUM	LOW
	Level B ~ "Probable" 1% – 10%	HIGH	HIGH	MEDIUM	LOW
	Level C ~ "Occasional" 0.1% – 1%	1/C	2/C	LOW	NEGLIGIBLE
	Level D ~ "Remote" 1-in-a-million – 0.1%	1/D	LOW	NEGLIGIBLE	NEGLIGIBLE
	Level E ~ "Improbable" < 10 ⁻⁶	LOW	NEGLIGIBLE	NEGLIGIBLE	NEGLIGIBLE

Figure 15: High-altitude Airdrop Risk Matrix

to be 1% of 1%, corresponding to a probability of death or severe injury of 10^{-4} . This is within the 90% confidence interval (9×10^{-5} to 2×10^{-3}) for the ‘actual’ probability based on the observed rate of 1 catastrophic event per 2000 person-flight hours (5×10^{-4}).³⁶

Based on the number of personnel flying and the number of high-altitude, unpressurized flights in the test program, the SRB concluded an overall probability of a catastrophic event during the program of 1/250 (0.4%).³⁷ This corresponds to a mishap severity-probability level of ***unlikely/catastrophic*** (1/C - HIGH RISK) per AFTCI 91-202 [12]. Based on these probability calculations, most SRB members felt constrained by the clear-cut numerical statistics from the ADRAC calculations. Other SRB members, particularly those with extensive experience flying operational high-altitude missions (e.g., aircrew and loadmaster), disputed a HIGH RISK characterization on the basis that “high-altitude airdrop is a routine operational mission.” Table 9 summarizes the various opinions of the SRB.

Table 9: C-17 high-altitude airdrop SRB judgments

Risk	severity level	mishap probability level	Comment
1/C	<i>catastrophic</i>	<i>unlikely</i> : 1-in-1000 to 1-in-100	$p = 1/250$ estimate
2/C	<i>critical</i>	<i>unlikely</i> : 1-in-1000 to 1-in-100	$BF > 10$
1/D	<i>catastrophic</i>	<i>highly unlikely</i> : 1-in-10 ⁶ to 1-in-1000	$\Theta^{(0)} < 10^{-3}$

From a Bayesian perspective, it is easy to understand and separate the viewpoints into three categories: those who agreed with the ***unlikely/catastrophic*** (1/C - HIGH RISK) categorization; those who doubted whether death was truly credible given the mitigations in place and favored

³⁶ We use a Gamma-Poisson conjugate prior (see Appendix B) to create a Bayesian estimate for the distribution of the rate of a rare event. With the Jeffreys (uninformed) prior $\sim \Gamma(0.5, 0)$ for the rate parameter λ , the posterior distribution becomes $\Gamma(1.5, 2000)$ after observing a single event in 2000 exposures (Figure B-3).

³⁷ Four flights with 10 personnel per flight yields 40 total exposures; the probability of at least one event in 40 trials is $1 - (1 - 10^{-4})^{40} = 0.004$.

an *unlikely/critical* (2/C - MED RISK) categorization; and those who, although acknowledging death as a credible consequence based on prior MDA mission experience, could not justify a subjective probability greater than 10^{-3} (*highly unlikely*) and favored maintaining the prior risk assessment of *catastrophic/highly unlikely* (1/D - MED RISK).

Differences of opinion or judgment among subject matter experts are not uncommon. Discounting deliberate ‘risk-hacking’ to ‘back-into’ a pre-determined risk category, the views and opinions of experts are usually sincere and genuinely held. Constructive conversations that achieve a consensus can be difficult or impossible if we lack a common framework or lexicon. Bayes Factors offer a solution. They provide a common lexicon by which to have meaningful discussions. Under this mental model for decision-making, differences of opinion are attributed to either differences in priors or differences in likelihood ratios. With the MDA test, those favoring a 1/D categorization had a different prior. Those who favored a 2/C categorization had a *BF* that strongly weighed the mitigations. Those who felt “bound by” the 1/C categorization based on the clear-cut mathematics were offered an opportunity for additional agency with a Bayesian approach that let members update their odds based on data.

We do not have to reach consensus in flight test; we have ultimate risk acceptance authorities whom we ask to make the final call. However, we need to provide our teams with the means to have meaningful and constructive discussions and to hear the diverse opinions that are key to realizing the wisdom of the crowd. Bayes Factors are a good framework for providing that structure.

Lesson Learned 12: *BFs* offer a framework with a common lexicon for discerning underlying reasons for differences in opinion and judgment (differences in posterior judgments reflect either a different prior or a different likelihood ratio)

V. Discussion & Lessons Learned

The examples in Section IV were deliberately selected to demonstrate the technique of using Bayes Factors to aid decision-making under uncertainty. Bayesian reasoning is a useful and reliable method for updating our understanding of an uncertain world. Bayes Factors are not a silver bullet. Probability will continue to be a measure of our subjective relationship with uncertainty. Bayes Factors do not change that, but they do provide a straightforward mechanism to work with our personal, subjective view of probability and uncertainty. Moreover, Bayes Factors allow us to sidestep the difficulty of estimating the probabilities of the observed data or evidence that is required in a direct application of Bayes’ Theorem. Systematic errors or erroneous assumptions about the underlying data are avoided because, as long as the same set of errors and assumptions is present in assessing the Bayes Factor, they cancel out. Even when the probabilities are clear, using Bayes Factors and odds ratios is frequently simpler than applying Bayes’ Theorem in conditional probability form (reference the examples in Sections IV.A and IV.D).

The posterior odds/probabilities are dependent on the prior odds/probabilities. The selection of prior probabilities is a persistent and sometimes contentious aspect of Bayesian analysis. Some critics are quick to dismiss the overall approach by claiming that priors are “just made up,”

thereby introducing arbitrary subjectivity into supposedly objective risk assessments. ET Jaynes³⁸ describes the problem: “An unfortunate impression has been created that rejection of personalistic probability automatically means the rejection of Bayesian methods in general. It will hopefully be shown here that this is not the case; the problem of achieving objectivity for prior probability assignments is not one of psychology or philosophy, but one of proper definitions and mathematical techniques, which is capable of rational analysis” [63].

In sharing the technique as described and demonstrated in this paper, the author has found the greatest consternation from ‘aspiring Bayesians’ to be determining the ‘right’ prior to use. The choice of priors, as well as the judgment of the relative strength of factors, is a judgment. To embrace a Bayesian approach to reasoning under uncertainty is to adopt the perspective that probabilities are explicit expressions of our current state of knowledge and degree of belief about uncertain things.

When we assign a prior probability for a test, and when we evaluate the strength of data about the system under test in the assessment of Bayes Factors, we are not arbitrarily picking numbers but systematically encoding our individual and collective engineering judgment. Flight test teams have the opportunity to incorporate and include experience with similar systems, similar tests, analysis of previous test results, understanding of team readiness, and an understanding of the physical principles of the test in the prior. Altogether, this offers a genuine assessment of what the team knows before conducting the test. The value in making explicit the assumptions and beliefs that would otherwise remain hidden is that they may be openly discussed by the team using a common framework and lexicon. Additionally, one of the mathematical properties of Bayesian inference ensures that with sufficient evidence, the choice of prior is irrelevant: different priors will converge to the same posterior distribution as data accumulates.

The real question is not whether we should use subjective priors, but whether we should make our inevitable subjectivity explicit and transparent. Without an explicit expression of our judgments, we leave them implicit and uninfluenced by the careful consideration of evidence. With an explicit Bayesian approach, our subjectivity is available for systematic revision. By exposing what we know or think we know, and by having a common reference point to compare our judgments with others, we gain a clear advantage from adopting a formal Bayesian approach in our flight test safety planning.

Differences of opinion should be expected and encouraged. With a Bayesian conceptual model for decision-making, the differences of opinion are either due to differences in priors or differences in likelihood ratios (*BFS*) based on the relative strength of the evidence. In revealing differences in prior odds or likelihood ratios, and in articulating our reasoning and judgment, we establish a framework for having meaningful and productive conversations that leverage the collective wisdom of the crowd. The author has also found these conversations informative, regardless of whether a consensus of opinions is reached, and helpful in shaping his ultimate judgment about a proposed test.

Many of the cognitive errors and biases described in Section II are compounded when

³⁸ Edwin Thompson Jaynes (1922–1998): American physicist and mathematician who fundamentally transformed probability theory and statistical inference; *Probability Theory: The Logic of Science* [32], is one of the most influential works in modern Bayesian statistics.

considering complex problems with complicated descriptions. Concise problem statements and decomposition of complex problems into simpler, more manageable, cognitively concise components can aid our analysis and judgment. By chunking data and evidence into separate pieces for piecewise consideration, Bayes Factors can help alleviate the overwhelming weight of cognitively digesting all the evidence at once. Additionally, by parsing our separate judgments and by making clear distinctions between priors and likelihoods, we can avoid some of the cognitive traps of decision-making. As an example of how the Bayes Factors approach may help avoid cognitive errors, consider the “prosecutor’s fallacy.”

The prosecutor’s fallacy is a critical cognitive error that confuses a statement about the likelihood ratio (the Bayes Factor) with a statement about the posterior probability given the evidence. *E.g.*, the prosecutor incorrectly switches $P(D_{DNA}|H_{Innocent})$ for $P(H_{Innocent}|D_{DNA})$. The fallacious prosecutor argues that since the probability that a DNA sample found at a crime scene would match a random person is 1-in-10-million, the probability that a defendant is innocent given a DNA sample match is 1-in-10-million. An equivalently erroneous argument would be swapping the statement “most popes are catholics” with “most catholics are popes.” The prosecutor’s fallacy is a form of base rate fallacy (Section IV.A) and has led to wrongful convictions and other miscarriages of justice.³⁹ In explicitly expressing our hypotheses, Bayes Factors, and odds ratios as we recommend here, we clarify otherwise complex statements and are less likely to make erroneous judgments due to confusing conditional probability statements.⁴⁰

Finally, because we are forced to consider alternative hypotheses when forming likelihood and odds ratios, we are less likely to be surprised by the unexpected event and may be less likely to fall victim to confirmation bias. *BFs* will not eliminate confirmation bias or other cognitive errors. But, by forcing us to be explicit about alternative hypotheses—and by assigning a numerical value to alternative outcomes—we have a better chance of avoiding the preferential tendency to see data that confirms our prior beliefs. A good guard against hubris when making decisions under uncertainty is to practice extreme humility and to be ready and willing to admit we may be wrong. Guarding against hubris is a healthy starting point for good judgment.

Does flight test really need another ‘new’ approach? Bayes Factors do not really represent a fundamentally different way of managing risk. *BFs* are merely a different way of organizing, communicating, and updating our state of knowledge and uncertainty. Having a common lexicon grounded in firm mathematical foundations is useful. As such, the 412 TW Test Safety Office, which teaches Test Safety planning at the USAF Test Pilot School (TPS), has adopted a Bayesian approach in its TPS curriculum. More work and education are necessary before Bayes Factors are widely adopted and routinely employed. The author has found the approach to be very useful in aiding decision-making under uncertainty. If the flight test community does as well, the benefits of a widespread adoption of common lexicon would be clear and likely endure.

³⁹ The O.J. Simpson and Sally Clark murder trials both include good, cautionary examples of the defense fallacy and the prosecutor’s fallacy [64, 65].

⁴⁰ Another example of the prosecutor’s fallacy and the use *BFs* to avoid an incorrect conclusion: consider a doping test which is 95% accurate, *i.e.*, $P(+\text{test}|\text{doper}) = 0.95$ and $P(-\text{test}|\text{clean}) = 0.95$. This yields a $BF = \frac{0.95}{0.05} = 19$. If the prevalence (proportion) of dopers in a sport is 1/50, then the posterior probability that someone is a doper given a positive test is only $(1/49)(19) = (19/49) \approx 28\%$, not 95% (the erroneous claim via the prosecutor’s fallacy).

A. Collected Lessons Learned on the Practical Use of Bayes Factors (*BFs*):

1. *BFs* are often simpler to apply than direct application of Bayes' Theorem, producing the "right" answer without needing to assess the problematic marginal probability of observed data (*BFs* are insensitive to assumptions about the data or evidence and guard against cognitive pitfalls associated with conditional probabilities)
2. Our judgments and conclusions are strongly influenced by our prior expectations
3. The *BF*-framework requires us to acknowledge alternate hypotheses explicitly (and may guard against confirmation bias, blind spots, and other cognitive errors)
4. *BFs* can mix both subjective and objective judgments
5. The *BF*-framework decomposes complex data sets into separate chunks and reveals the relative strength of different pieces of evidence (combining multiple pieces of evidence is as simple as multiplying *BFs*)
6. *BFs* provide a systematic structure for meaningful conversations and communicating our experience, judgment, and decisions
7. *BFs* can account for and incorporate other factors that were not part of the original problem statement
8. *BFs* are straightforward to perform real-time and via "back-of-the-envelope" calculations
9. *BFs* provide a common framework to express and expose the full range of subject matter expert opinions; the range of opinions provides a good estimate for uncertainty and for bounding judgments with confidence intervals
10. *BFs* are also applicable and useful for making technical, programmatic, schedule, and other risk assessments
11. *BFs* are useful for tracking changing judgments over time and for communicating those judgments with stakeholders; includes communicating progress towards an objective and the relative importance of factors for "getting to yes" (avoiding a sequence of "bring me a rock" drills)
12. *BFs* offer a framework with a common lexicon for discerning underlying reasons for differences in opinion and judgment (differences in posterior judgments reflect either a different prior or a different likelihood ratio)

VI. Conclusion

*There is no such thing as absolute certainty, but there is
assurance sufficient for the purposes of human life*

John Stuart Mill [66]

Risk management, in flight test or otherwise, is about exercising judgment and making decisions under uncertainty. All risk management frameworks and decision-making models incorporate an explicit or implicit assessment of the likelihood of success or failure of a particular outcome. Accurate probability assessments for many practical problems of epistemic uncertainty in the “real-world” are difficult. Even the best engineering judgments sometimes vary by several orders of magnitude.

Guided by the twin observations that 1) objective probability does not exist and 2) the wisdom of the crowd is effective, we embraced Bayes Factors as a means of having a common framework and lexicon to improve our judgment of risk and decision-making under uncertainty. Bayes Factors provided a means to practice Bayesian reasoning, systematically incorporating new data and evidence, while sidestepping some of the challenges of assessing the marginal probability of the data.

In the use of Bayes Factors, we found value in being forced to state our priors, revealing our assumptions explicitly. By explicitly admitting and considering alternative hypotheses, we expect to benefit from a decreased likelihood of being surprised by the unexpected. We found clarity in discussing and acknowledging our mental models when describing likelihood ratios. Bayes Factors will not altogether eliminate cognitive biases such as confirmation bias, but they do help guard against them.

We should never forget that uncertainty will always lie at the heart of flight test. Our goal is to be “*intelligent and full of doubt*”—to be consciously aware of our ignorance while making the best decisions under uncertainty that we can, and to have the extreme humility to admit that we could be wrong. Becoming Bayesian is a good guard against hubris and being cocksure. Absolute certainty may not exist, but *Becoming Bayesian* may grant sufficient assurance for human life.

Acknowledgments & Disclaimer

I have spent two decades thinking about uncertainty. As a student of risk management, I am indebted to many in our flight test community for the discussions, ideas, feedback, critiques, and dedication to making our profession as safe and rigorous as possible. Since learning about the Bayes Factor approach while serving as a Permanent Professor at USAFA, I have significantly benefited from multiple, extended intellectual exchanges on the topic. These helped inform and clarify my understanding, practice, and explanation. Additionally, I am grateful for the patience and indulgence of the men and women of the 412th Test Wing as we practiced and implemented the application of Bayes Factors in our flight test risk management at Edwards AFB.

It is challenging to fully credit everyone who has helped to hone the ideas in this paper.

Many have contributed in many different ways. In particular, I am grateful for the discussions and feedback from Maj Gen Scott “Nova” Cain, Brig Gen (ret) Matthew Higer, Mr Art Huber, Mr Tim “Batman” West, Col Matthew “KITT” Caspers, Col (ret) Andy “Sonic” Freeborn, Col James “Nut” Gresham, Col James “T-Pain” Hayes, Col Dan “Animal” Javorsek, Col (ret) Mark “SCIPR” Jones Jr., Col Maryann “Yetti” Karlen, Col Brian “Bandit” Neff, Col Ryan “Hulk” Sanford, Col James “Fangs” Valpiani, Mr Jesus Arzate, Dr James Brownlow, Mr Nathan “CAP’N” Cook, Mr Kirk Harwood, Mr Richard Jones, Mr Wei “FUG” Lee, Mr Chris “NORAD” Liebmann, Mr Dan “Dan-O” Osburn, Mr Dave Vanhoy, Lt Col (ret) Tyler “Matrix” Robarge, Lt Col Timothy “Speed” Lau, Lt Col Dan Edelstein, Lt Col Matthew “Fear” Gray, Lt Col Robert “Bear” Newton, Lt Col Carlos “Yardman” Pinedo, Dr Robert Poulson, and Lt Col (ret) Donald “Whiz” Sheesley.

The views expressed are those of the author and do not reflect the official policy or position of the U.S. Air Force, the Department of Defense, or the U.S. Government. This material has been approved for public release and unlimited distribution.

References

- [1] Bertrand Russell. “The Triumph of Stupidity”. In: *Mortals and Others: Bertrand Russell’s American Essays, 1931-1935*. Ed. by Harry Ruja. Routledge, pp. 27–28.
- [2] Douglas P Wickert. “Risk Awareness: A New Framework for Risk Management in Flight Test”. In: *Proceedings to the SETP 62nd Annual Symposium* (2018). URL: https://github.com/dpwickert/dpwickert.github.io/blob/aa0db2414e51263616b9cda6edb68a9ab07f95f2/pdfs/Risk_Awareness_SETP_2018.pdf.
- [3] Mica R Endsley. “Design and evaluation for situation awareness enhancement”. In: *Proceedings of the Human Factors Society annual meeting*. Vol. 32. 2. Sage Publications Sage CA: Los Angeles, CA. 1988, pp. 97–101.
- [4] U.S. Department of Defense. *MIL-STD-882E w/CHANGE 1, Department of Defense Standard Practice for System Safety*. Office of the Under Secretary of Defense for Acquisition and Sustainment, July 2023. URL: https://www.cto.mil/wp-content/uploads/2025/07/MIL-STD-882E-w_CHANGE-1.pdf.
- [5] D.J. Spiegelhalter. *The Art of Uncertainty: How to Navigate Chance, Ignorance, Risk and Luck*. W.W. Norton & Company, 2025. ISBN: 9780241658635.
- [6] Amos Tversky and Daniel Kahneman. “Judgment under Uncertainty: Heuristics and Biases”. In: *Science* 185.4157 (1974), pp. 1124–1131. DOI: 10.1126/science.185.4157.1124.
- [7] Gerd Gigerenzer, Peter M. Todd, and ABC Research Group. *Simple heuristics that make us smart*. New York, USA: Oxford University Press, 1999.
- [8] Herbert A. Simon. *Administrative Behavior: A Study of Decision-Making Processes in Administrative Organization*. 4th ed. New York: Free Press, 1997.
- [9] James Clerk Maxwell. “Book review: Watson’s Kinetic Theory of Gases”. In: *Nature* 16 (1877), pp. 242–246.
- [10] International Organization for Standardization. *ISO 31000:2018, Risk management—Guidelines*. Standard. Geneva, Switzerland: International Organization for Standardization, 2018. URL: <https://www.iso.org/standard/65694.html>.
- [11] Department of the Air Force. *Risk Management Guidelines and Tools*. Tech. rep. DAFFPAM 90-803. Department of the Air Force, Mar. 2022.
- [12] Air Force Test Center. *Test Safety Program*. AFTCI 91-202. U.S. Air Force. Nov. 2022. URL: <https://static.e-publishing.af.mil/production/1/aftc/publication/aftci91-202/aftci91-202.pdf>.

- [13] Department of the Air Force. *DAFI 90-802, Risk Management*. Tech. rep. Published on the Air Force e-Publishing website. Headquarters, Air Force Safety, Sept. 2024. URL: https://static.e-publishing.af.mil/production/1/af_se/publication/afi90-802/dafi90-802.pdf.
- [14] Leveson, Nancy G. and Thomas, John P. *STPA handbook*. MIT Partnership for Systems Approaches to Safety and Security (PSASS). Cambridge, Massachusetts, U.S., 2018.
- [15] Michael Stamatelatos and Homayoon Dezfuli. *Probabilistic Risk Assessment Procedures Guide for NASA Managers and Practitioners*. Tech. rep. NASA/SP-2011-3421. Washington, DC: NASA Headquarters, 2011.
- [16] E. T. Bell. *The Development of Mathematics*. 2nd. New York: McGraw-Hill Book Company, 1945.
- [17] Pierre-Simon Laplace. *Théorie analytique des probabilités*. Paris: Courcier, 1814.
- [18] Jacob Bernoulli. *Ars Conjectandi*. Basel: Thurnisius, 1713.
- [19] John Maynard Keynes. *A Treatise on Probability*. London: Macmillan & co., 1921.
- [20] Rudolf Carnap. *Logical Foundations of Probability*. Chicago: University of Chicago Press, 1950.
- [21] Thomas Bayes. “An Essay towards solving a Problem in the Doctrine of Chances”. In: *Philosophical Transactions of the Royal Society of London* 53 (1763), pp. 370–418.
- [22] Bruno de Finetti. *Theory of Probability: A Critical Introductory Treatment*. 2 vols., English translation of *Teoria delle Probabilità* (1937). London: John Wiley & Sons, 1974.
- [23] Richard von Mises. *Wahrscheinlichkeit, Statistik und Wahrheit*. English translation: *Probability, Statistics and Truth*, Dover, 1957. Vienna: Julius Springer, 1928.
- [24] Ronald A. Fisher. *On the Mathematical Foundations of Theoretical Statistics*. Vol. 222. 594–604. 1922, pp. 309–368.
- [25] Jerzy Neyman and E. S. Pearson. “On the Problem of the Most Efficient Tests of Statistical Hypotheses”. In: *Philosophical Transactions of the Royal Society of London. Series A* 231 (1933), pp. 289–337. DOI: [10.1098/rsta.1933.0009](https://doi.org/10.1098/rsta.1933.0009).
- [26] Karl R. Popper. *The Propensity Interpretation of Probability*. London: The British Academy, 1959.
- [27] Andrey N. Kolmogorov. *Grundbegriffe der Wahrscheinlichkeitsrechnung*. English translation: *Foundations of the Theory of Probability*, Chelsea Publishing, 1950. Berlin: Springer, 1933.
- [28] Alan Hájek. “Interpretations of Probability”. In: *The Stanford Encyclopedia of Philosophy*. Ed. by Edward N. Zalta and Uri Nodelman. Winter 2023. Accessed on 29 Aug 2025. Metaphysics Research Lab, Stanford University, 2023.
- [29] Frank P Ramsey. “Truth and probability”. In: *Readings in formal epistemology: Sourcebook*. Springer, 1926, pp. 21–45.
- [30] A. Clayton. *Bernoulli's Fallacy: Statistical Illogic and the Crisis of Modern Science*. Columbia University Press, 2021. ISBN: 9780231553353.
- [31] H. Jeffreys. *The Theory of Probability*. Oxford Classic Texts in the Physical Sciences. OUP Oxford, 1998. ISBN: 9780191589676. URL: <https://books.google.com/books?id=vh9Act9rtzQC>.
- [32] E.T. Jaynes and G.L. Bretthorst. *Probability Theory: The Logic of Science*. Cambridge University Press, 2003. ISBN: 9780521592710.
- [33] N. Silver. *The Signal and the Noise: Why So Many Predictions Fail—but Some Don't*. Penguin Publishing Group, 2012. ISBN: 9781101595954.
- [34] A. Gelman et al. *Bayesian Data Analysis*. Chapman & Hall/CRC Texts in Statistical Science. CRC Press, 2013. ISBN: 9781439898208.
- [35] C. Perrow. *Normal Accidents: Living with High Risk Technologies*. Princeton paperbacks. Princeton University Press, 1999. ISBN: 9780691004129.
- [36] David L. Brumley. ”*Making Safety Happen” Through Probabilistic Risk Assessment at NASA*. Tech. rep. 20200001592. National Aeronautics and Space Administration (NASA), Mar. 2020. URL: <https://ntrs.nasa.gov/citations/20200001592>.

- [37] Teri Hamlin et al. *SHUTTLE RISK PROGRESSION BY FLIGHT*. Presentation at NASA. Presented on June 24, 2011. NASA, Johnson Space Center; SAIC; MSFC-BTI, June 2011.
- [38] F. Galton. "Vox populi". In: *Nature* 75.1949 (1907), p. 7.
- [39] P.E. Tetlock and D. Gardner. *Superforecasting: The Art and Science of Prediction*. Crown, 2015. ISBN: 9780804136709.
- [40] D. Gardner. *Future Babble: Why Pundits Are Hedgehogs and Foxes Know Best*. Penguin Publishing Group, 2012. ISBN: 9780452297579. URL: <https://books.google.com/books?id=6qNPEAAAQBAJ>.
- [41] I. Berlin, H. Hardy, and M. Ignatieff. *The Hedgehog and the Fox: An Essay on Tolstoy's View of History - Second Edition*. Princeton University Press, 2013. ISBN: 9780691156002. URL: <https://books.google.com/books?id=V2iYDwAAQBAJ>.
- [42] J. Surowiecki. *The Wisdom of Crowds*. Knopf Doubleday Publishing Group, 2005. ISBN: 9780307275059.
- [43] Mohamad Hjeij and Arnis Vilks. "A brief history of heuristics: how did research on heuristics evolve?" In: *Humanities and Social Sciences Communications* 10.1 (Mar. 2023). DOI: 10.1057/s41599-023-01542-z. URL: <https://doi.org/10.1057/s41599-023-01542-z>.
- [44] D. Kahneman. *Thinking, Fast and Slow*. Farrar, Straus and Giroux, 2011. ISBN: 9781429969352.
- [45] G. Gigerenzer. *Calculated Risks: How to Know When Numbers Deceive You*. Calculated Risks: How to Know when Numbers Deceive You. Simon & Schuster, 2002. ISBN: 9780743254236.
- [46] F.H. Knight. *Risk, Uncertainty and Profit*. Hart, Schaffner and Marx prize essays. Houghton Mifflin, 1921. ISBN: 9781548563509.
- [47] Andy Stirling. "Keep it complex". In: *Nature* 468.7327 (2010), pp. 1029–1031.
- [48] Frank Plumpton Ramsey. "Probability and partial belief". In: (1961).
- [49] Herbert A Simon. "A Behavioral Model of Rational Choice". In: *Quarterly Journal of Economics* 69.1 (1955), pp. 99–118.
- [50] Gerd Gigerenzer and Daniel G. Goldstein. "Reasoning the fast and frugal way: Models of bounded rationality". In: *Psychological Review* 103 (1996), pp. 650–669. DOI: 10.1037/0033-295X.103.4.650.
- [51] G. Gigerenzer. *Risk Savvy: How to Make Good Decisions*. Penguin Publishing Group, 2014. ISBN: 9780698151437.
- [52] G.A. Klein. *Sources of Power: How People Make Decisions*. The MIT Press. MIT Press, 1999. ISBN: 9780262260862.
- [53] Jacob Q. Robinson. "Intuitive Judgment and Strategic Decisions". DTIC Accession Number: AD1160094. Master of Military Art and Science Thesis. Fort Leavenworth, Kansas: Advanced Strategic Leadership Studies Program (ASLSP), Sept. 2020.
- [54] Nathan Cook. *Personal communication*. correspondence with the author. Sept. 2025.
- [55] Hugo Bottemanne. "Bayesian brain theory: Computational neuroscience of belief". In: *Neuroscience* 566 (2025), pp. 198–204. ISSN: 0306-4522. DOI: <https://doi.org/10.1016/j.neuroscience.2024.12.003>. URL: <https://www.sciencedirect.com/science/article/pii/S0306452224007048>.
- [56] Robert Bain. "Are our brains Bayesian?" In: *Significance* 13.4 (2016), pp. 14–19.
- [57] Robert E. Kass and Adrian E. Raftery. "Bayes Factors". In: *Journal of the American Statistical Association* 90.430 (1995), pp. 773–795. ISSN: 01621459, 1537274X. URL: <http://www.jstor.org/stable/2291091> (visited on 09/09/2025).
- [58] Amos Tversky and Daniel Kahneman. *Evidential Impact of Base Rates*. Technical Report TR-4. Supported by the Office of Naval Research under Contract N00014-79-C-0077, Work Unit NR 197-058. Stanford, CA: Stanford University, Department of Psychology, May 1981. DOI: 10.21236/ADA099501. URL: <https://apps.dtic.mil/sti/html/tr/ADA099501/index.html>.
- [59] Martin Gardner. "Mathematical Games". In: *The Second Scientific American Book of Mathematical Puzzles and Diversions*. Reprinted in. Oct. 1959, pp. 180–182.
- [60] Steve Selvin. "A problem in probability (letter to the editor)". In: *The American Statistician* 29.1 (1975), pp. 67–71. DOI: 10.1080/00031305.1975.10479121.

- [61] Marilyn vos Savant. “Ask Marilyn”. In: *Parade Magazine* (Sept. 1990). Question posed in a letter by Craig F. Whitaker., p. 16.
- [62] Marilyn vos Savant. “Ask Marilyn”. In: *Parade Magazine* (Dec. 1990), p. 25.
- [63] Edwin T. Jaynes. “Prior Probabilities”. In: *IEEE Transactions on Systems Science and Cybernetics* SSC-4.3 (1968), pp. 227–241. DOI: 10.1109/TSSC.1968.300117.
- [64] I. J. Good. “When batterer turns murderer”. In: *Nature* 375 (1995), p. 541. DOI: 10.1038/375541a0.
- [65] Royal Statistical Society. *Royal Statistical Society concerned by issues raised in Sally Clark case (News Release)*. Royal Statistical Society. Statement issued October 23, 2001, regarding the misuse of statistics in the courts. Oct. 2001. URL: <https://rss.org.uk/RSS/media/File-library/Membership/Sections/2020/Sally-Clark-RSS-statement-2001.pdf>.
- [66] John Stuart Mill. *On Liberty*. John W. Parker and Son, 1859.
- [67] Grant Sanderson. *Bayes theorem, the geometry of changing beliefs*. Dec. 2019. URL: <https://youtu.be/HZGCoVF3YvM> (visited on 07/27/2025).
- [68] Hoesung Lee et al. “Climate change 2023 synthesis report summary for policymakers”. In: *CLIMATE CHANGE 2023 Synthesis Report: Summary for Policymakers* (2024).
- [69] Rebecca Holliday et al. “Record-breaking May heat in the UK: contrasting the extreme temperatures of 2024 and 1944 using climate attribution”. In: *Weather* (2024).

Appendix A: Derivation of Bayes Factor

The Bayes Factor is the likelihood ratio that assesses the relative weight of evidence, data, or judgment between two competing hypotheses. Sir Harold Jeffreys first introduced the concept in his 1939 book, *Theory of Probability* [31]. A more widespread adoption by statisticians was encouraged and emphasized by Kass and Raftery in 1995 [57].

We derive the Bayes Factor from Bayes' Theorem and describe its interpretation. We will see that the Bayes Factor emerges naturally from Bayes' theorem when we compare competing hypotheses. Noting the symmetry of probability between two states, A and B , we have

$$P(A \cap B) = P(A|B)P(B) = P(B \cap A) = P(B|A)P(A) \quad (\text{A-1})$$

where $P(A|B)$ is the probability of A given B . Bayes' theorem naturally follows:

$$P(A|B) = \frac{P(B|A)P(A)}{P(B)} \quad (\text{A-2})$$

The relationship is depicted geometrically in Figure A-1.⁴¹

Expressing (A-2) in terms of the probability of observing a hypothesis, H , given some data or evidence D , we obtain Bayes' theorem in its standard form:

$$P(H|D) = \frac{P(D|H)P(H)}{P(D)} = \frac{P(D|H)P(H)}{P(H)P(D|H) + P(H^c)P(D|H^c)} \quad (\text{A-3})$$

where:

- $P(H|D)$ is the *posterior probability* of hypothesis H given data D
- $P(D|H)$ is the *likelihood* of observing data D given hypothesis H
- $P(H)$ is the *prior probability* of hypothesis H
- $P(D)$ is the *marginal likelihood* or evidence, which normalizes the probability to ensure it lies between 0 and 1 (this will prove problematic)

The marginal likelihood of observing the data, $P(D)$, presents a computational challenge because it requires integrating or summing over all possible hypotheses. However, we can elegantly sidestep this issue by considering the alternative hypothesis and examining the ratio of posterior probabilities.

Let H^c represent the complement of hypothesis H . We can also write Bayes' theorem for this complementary hypothesis:

$$P(H^c|D) = \frac{P(D|H^c)P(H^c)}{P(D)} \quad (\text{A-4})$$

⁴¹ Inspired by Grant Sanderson's 3Blue1Brown video [67]

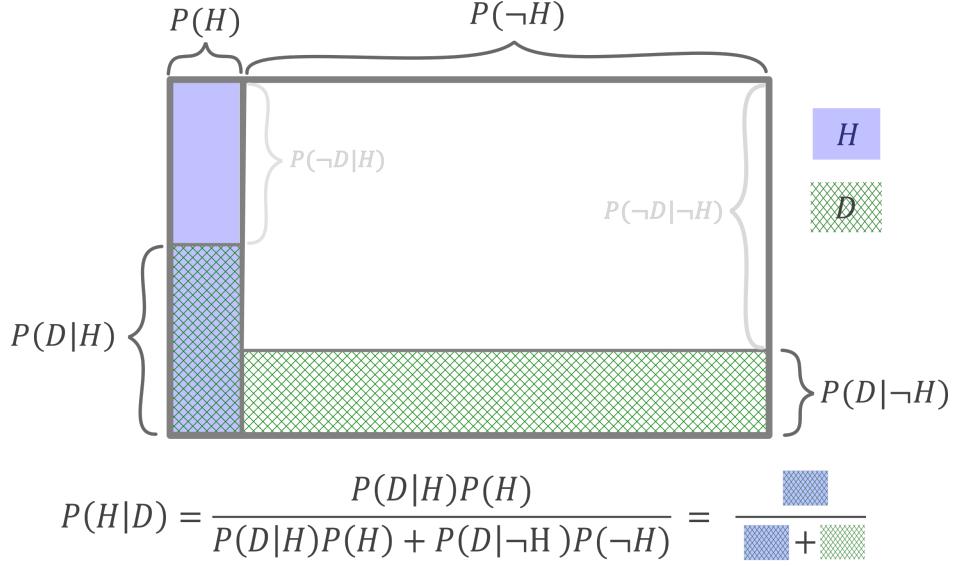


Figure A-1: Geometric Interpretation of Bayes' Theorem

The ratio of the two posterior probabilities, (A-3) and (A-4), is

$$\frac{P(H|D)}{P(H^c|D)} = \frac{\frac{P(D|H)P(H)}{P(D)}}{\frac{P(D|H^c)P(H^c)}{P(D)}} \quad (\text{A-5})$$

Note that the problematic term assessing the probability of observing the data, $P(D)$, in (A-3)

$$P(D) = P(H)P(D|H) + P(H^c)P(D|H^c)$$

appears in both the numerator and denominator in (A-5), allowing us to cancel it. Thus,

$$\frac{P(H|D)}{P(H^c|D)} = \frac{P(D|H)P(H)}{P(D|H^c)P(H^c)} \quad (\text{A-6})$$

The left-hand side is the **posterior odds** for hypothesis H given the data D . The elegance of this formulation is that we have eliminated the troublesome marginal likelihood $P(D)$ entirely. Rewriting (A-6) to separate the likelihood ratio from the prior odds yields

$$\frac{P(H|D)}{P(H^c|D)} = \frac{P(D|H)}{P(D|H^c)} \cdot \frac{P(H)}{P(H^c)} \quad (\text{A-7})$$

To use Bayes Factors we have to express our probability in terms of odds ratios, Θ . A probability, p , is converted to odds by

$$\Theta = \frac{p}{1-p} \quad (\text{A-8})$$

Odds are converted back into probabilities by

$$p = \frac{\Theta}{\Theta + 1} \quad (\text{A-9})$$

Defining the **prior odds** as

$$\Theta^{(0)} = \frac{P(H)}{P(H^c)} \quad (\text{A-10})$$

and the **posterior odds**

$$\Theta^{(1)} = \frac{P(H|D)}{P(H^c|D)} \quad (\text{A-11})$$

yields the definition of the **Bayes Factor**:

$$BF = \frac{P(D|H)}{P(D|H^c)} \quad (\text{A-12})$$

Equation A-6 can now be rewritten in final form using the Bayes Factor:

$$\Theta^{(1)} = BF \cdot \Theta^{(0)} \quad (\text{A-13})$$

Interpretation

The Bayes Factor, Equation A-12, represents the strength of evidence provided by the data D in favor of hypothesis H over hypothesis H^c . It quantifies how many times more likely the observed data is under hypothesis H compared to hypothesis H^c .

Interpreting the strength of the Bayes Factor:

- $BF \gg 1$: Evidence provides strong support for hypothesis H
- $BF > 1$: Evidence favors hypothesis H
- $BF \approx 1$: Evidence is equally consistent with both hypotheses
- $BF < 1$: Evidence favors hypothesis H^c
- $BF \ll 1$: Evidence provides strong support for hypothesis H^c

Combined Bayes Factor for N pieces of data is the product of the individual BF_i 's so that

$$BF = \prod_{i=1}^N BF_i \quad (\text{A-14})$$

The recent use of attribution methods in climate science [68, 69], which involves comparing the likelihood of particular climatic events using Monte Carlo simulations of two different models, with and without climate change effects, is an interesting example of the use of Bayes Factors.

Appendix B: Conjugate Prior Distributions in Bayesian Reasoning

A conjugate prior⁴² is a prior distribution that, when combined with a particular likelihood function through Bayes' theorem, produces a posterior distribution from the same family as the prior with new parameters. This mathematical convenience allows us to compute our posterior analytically rather than resorting to complex numerical methods. This appendix summarizes two conjugate priors: the Beta and Gamma distributions.

Both the Beta and Gamma distributions are readily available in statistical software packages, making them straightforward to use.

- **Excel:** use the `BETA.DIST(x, α, β , cumulative)` and `GAMMA.DIST(x, α, β , cumulative)` functions; the boolean *cumulative* parameter returns either the CDF (TRUE) or PDF (FALSE)
- **SciPy:** Python's `scipy.stats` library contains `beta.pdf(x, α, β)`, `beta.cdf(x, α, β)`, `gamma.pdf(x, α, β)`, and `gamma.cdf(x, α, β)`
- **NumPy:** Python's NumPy library contains `np.random.beta(α, β , size)` and `np.random.gamma(α, β , size)` for random sampling
- **MATLAB:** `betapdf(x, α, β)`, `betacdf(x, α, β)`, `gampdf(x, α, β)`, and `gamcdf(x, α, β)` with corresponding random generators `betarnd(α, β)` and `gamrnd(α, β)`

Note that software packages may use different parameterizations, particularly for the Gamma distribution's scale vs. rate parameter.

Both the Beta and Gamma distributions are used as conjugate priors in Examples in Section IV.

B-I. Beta Distribution

The Beta distribution is a continuous probability distribution defined on the interval $[0, 1]$, which is a convenient conjugate prior for modeling proportions or Bernoulli/binomial trials in Bayesian statistics. The Beta distribution is parameterized by two positive parameters ($\alpha, \beta > 0$) which control the shape of the distribution. The probability density function is

$$f(x; \alpha, \beta) = \frac{x^{\alpha-1}(1-x)^{\beta-1}}{\mathcal{B}(\alpha, \beta)} \quad (\text{B-1})$$

where $\mathcal{B}(\alpha, \beta)$ is the Beta function which serves as a normalization constant. The Beta function is defined as

$$\mathcal{B}(\alpha, \beta) = \frac{\Gamma(\alpha)\Gamma(\beta)}{\Gamma(\alpha + \beta)} \quad (\text{B-2})$$

where $\Gamma(\cdot)$ is the Gamma function:

$$\Gamma(z) = \int_0^\infty t^{z-1} e^{-t} dt \quad (\text{B-3})$$

⁴² Conjugate comes from Latin *conjugatus*, meaning “joined” or “yoked together.”

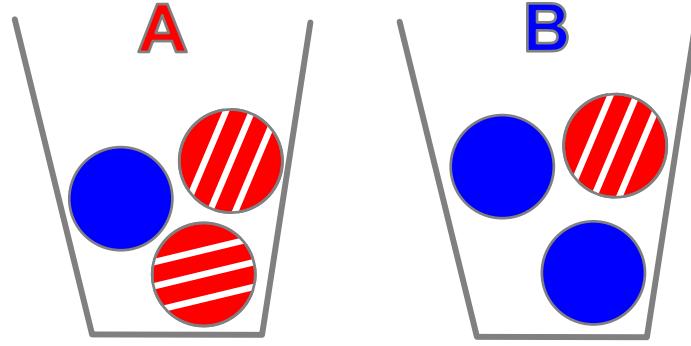


Figure B-1: Two Bags with Three Balls Example

The Gamma function satisfies the recurrence relation $\Gamma(z + 1) = z\Gamma(z)$, so that it is often known as the ‘factorial function,’ owing to the property that

$$\Gamma(n + 1) = n! \quad (\text{B-4})$$

Table B-1: Statistics for the Beta Distribution ($\mathcal{B}(\alpha, \beta)$)

statistic	value	support
mean	$\mu = \frac{\alpha}{\alpha+\beta}$	$\alpha, \beta > 0$
mode	$\frac{\alpha-1}{\alpha+\beta-2}$	$\alpha, \beta > 1$
variance	$\sigma^2 = \frac{\alpha\beta}{(\alpha+\beta)^2(\alpha+\beta+1)}$	

The interpretation of the shape parameters is intuitive: α is the number of ‘successes’ and β represents the number of ‘failures’ in a series of Bernoulli trials (a binomial experiment). With a $\mathcal{B}(\alpha_{prior}, \beta_{prior})$ distribution, and x successes in n additional trials, the posterior distribution is also Beta-distributed and given by

$$\mathcal{B}(\alpha_{posterior}, \beta_{posterior}) = \mathcal{B}(\alpha_{prior} + x, \beta_{prior} + n - x) \quad (\text{B-5})$$

Example: two bags with three colored balls

To illustrate the application of Beta distributions in Bayesian reasoning, consider the contrived academic example in which we have two bags, each containing three balls with two different colors:

- Bag A: 2 red balls, 1 blue ball (probability of drawing blue = $\frac{1}{3}$)
- Bag B: 1 red ball, 2 blue balls (probability of drawing blue = $\frac{2}{3}$)

One of the bags is selected at random, and we then draw balls from that bag (with replacement) recording the colors. We use Bayesian reasoning to update our prior belief about which bag was

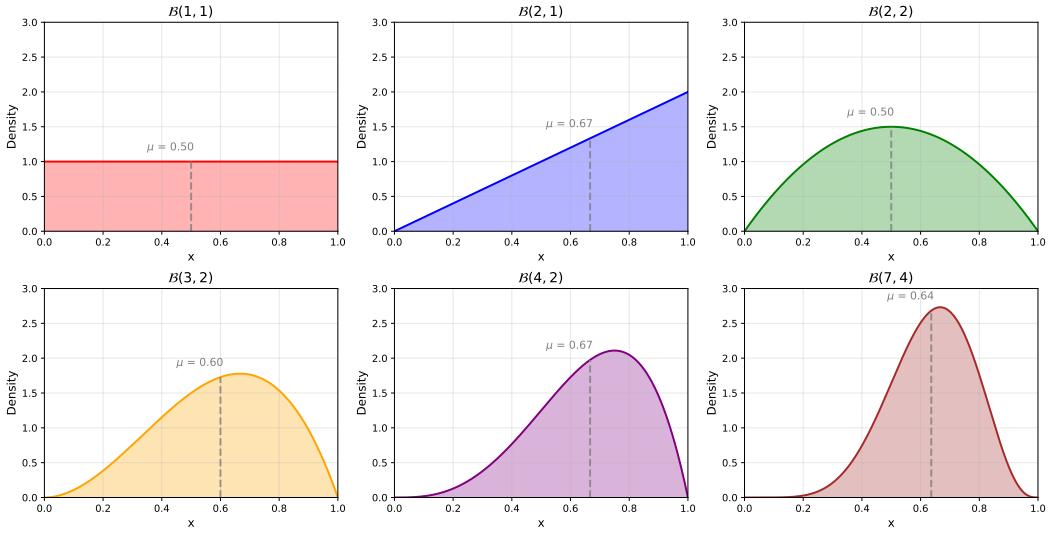


Figure B-2: Posterior Beta Distributions for the Probability of drawing a Blue ball

selected and what the probability of drawing a blue ball is using the Beta distribution as our conjugate prior.

Since a bag is selected at random, we have even prior odds for selecting Bag B with two blue balls, thus $\Theta_B^{(0)} = 1 : 1$. Let's take the sequence of draws (with replacement) given in Table B-2.

Table B-2: Two-Bag Example

	draws	BF_i^{\dagger}	$\Theta_B^{(i)}$	$p(B D_i)$	\mathcal{B}_i
0	—	—	$\Theta_B^{(0)} = 1 : 1$	$p = 1/2$	$\mathcal{B}_0(1, 1)$
1	B	$BF_1 = \frac{2/3}{1/3} = 2$	$\Theta_B^{(1)} = 2 : 1$	$p = 2/3$	$\mathcal{B}_1(2, 1)$
2	B,R	$BF_2 = \frac{1/3}{2/3} = 1/2$	$\Theta_B^{(2)} = 1 : 1$	$p = 1/2$	$\mathcal{B}_2(2, 2)$
3	B,R,B	$BF_3 = \frac{2/3}{1/3} = 2$	$\Theta_B^{(3)} = 2 : 1$	$p = 2/3$	$\mathcal{B}_3(3, 2)$
4	B,R,B,B	$BF_4 = \frac{2/3}{1/3} = 2$	$\Theta_B^{(4)} = 4 : 1$	$p = 4/5$	$\mathcal{B}_4(4, 2)$
	⋮				
9	B,R,B,B,R,R,B,B,B	—	$\Theta_B^{(9)} = 8 : 1$	$p = 8/9$	$\mathcal{B}_9(7, 4)$

[†] BF_i is the Bayes Factor for the draw of the i^{th} ball alone

After a sequence of 4 draws consisting of 3 blue balls and 1 red ball, we are 80% confident that we had originally selected Bag B. The posterior distribution for the probability of drawing a blue ball is given in Figure B-2 for each of the six cases in Table B-2.

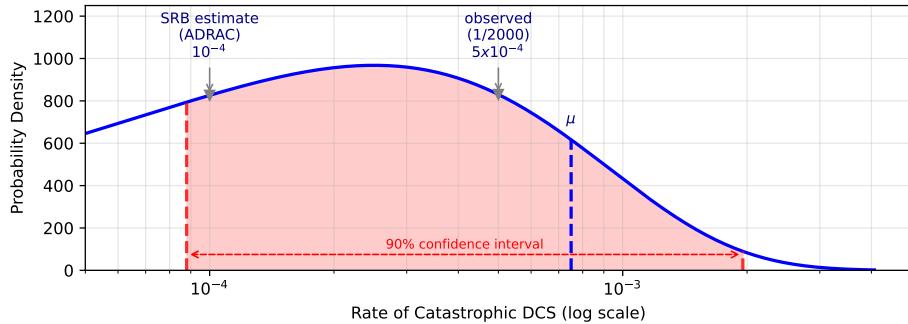


Figure B-3: $\Gamma(1.5, 2000)$ Posterior Distribution for Expected Rate of a Rare Event

B-II. Gamma Distribution

The Gamma distribution is a continuous probability distribution defined on the interval $(0, \infty)$ that serves as a conjugate prior for modeling positive-valued parameters such as rates, scales, and precision parameters. The exponential, Erlang, and chi-squared distributions are special cases of the Gamma distribution. The Gamma distribution is parameterized by two positive parameters: shape parameter $\alpha > 0$ and a rate parameter $\lambda > 0$. The probability density function is

$$f(x; \alpha, \lambda) = \frac{\lambda^\alpha}{\Gamma(\alpha)} x^{\alpha-1} e^{-\lambda x} \quad (\text{B-6})$$

where $\Gamma(\alpha)$ is the Gamma function defined by (B-3). An alternative parameterization uses the scale parameter $\theta = 1/\lambda$, yielding

$$f(x; \alpha, \theta) = \frac{1}{\Gamma(\alpha)\theta^\alpha} x^{\alpha-1} e^{-x/\theta} \quad (\text{B-7})$$

Table B-3: Statistics for the Gamma Distribution, $\Gamma(\alpha, \lambda)$

	statistic	value	support
mean	$\mu = \frac{\alpha}{\lambda}$	$\alpha, \lambda > 0$	
mode	$\frac{\alpha-1}{\lambda}$	$\alpha > 1$	
variance	$\sigma^2 = \frac{\alpha}{\lambda^2}$		

The Gamma distribution exhibits important conjugacy properties with several likelihood functions. For Poisson data with rate parameter λ , as was used in Section IV.H, a $\Gamma(\alpha_{prior}, \lambda_{prior})$ prior combined with observed count data $\sum_{i=1}^n x_i$ over n observations yields the posterior

$$\Gamma(\alpha_{posterior}, \lambda_{posterior}) = \Gamma(\alpha_{prior} + \sum_{i=1}^n x_i, \lambda_{prior} + n) \quad (\text{B-8})$$

For exponential data with rate parameter λ , the Gamma distribution serves as a conjugate prior, making it useful for modeling failure rates and other positive, continuous random variables.