

Utility Representation - Table of Contents

<i>The Average Man</i>	2
<i>...and Expected Value</i>	3
<i>Part I: Probability Measures</i>	4
<i>Probability: Two sides of a Coin</i>	4
<i>Small Worlds and Statistical Inference</i>	6
<i>Modeling: Improper Assumptions and Skewed Expectations</i>	8
<i>Frequentism: Inference from Expected Frequency</i>	10
<i>Bayesian Inference: Inference to Expected Value</i>	11
<i>An Example: Expected Website Returns</i>	12
<i>Part II: Utility Curves</i>	14
<i>The Stakes: From Utility to Indifference</i>	14
<i>Optimising Utility</i>	15
<i>Part III: Representation Theorems</i>	16
<i>Rational Preference: From Indifference to Utility</i>	16
<i>Von Neumann-Morgenstern Representation: Decision Under Risk</i>	18
<i>Bolker-Jeffrey Representation: Decision Under Uncertainty</i>	19
<i>Part IV: Machine Learning and the Individual</i>	23
<i>Customer Representation: A Segmentation Approach</i>	23
<i>PCA and Segmentation</i>	26
<i>Conclusion: Construct Criticism</i>	28

Utility Theory and the Representation of Preference

Nathaniel Forde

March 1, 2021

The Average Man

In the 1840s the Average man stalked the nightmares of Augustin Cournot. The mathematician was haunted by a melange of imagined parts. A Frankenstein's monster of mismatched limbs, variously soldered, stapled or sown at the joints. He worried that no one but a "physical monstrosity" could exhibit the average, weight, height and other mean attributes in one body. But fantastical fears were no bar to a useful mathematical fiction. The notion of averaging was a technological breakthrough - applications of averaging multiplied without cease: polling, gambling, forecasting - statistics were recorded everywhere, compounding one on another; averages of averages tenuously tethered to observable facts by layers of abstraction and scales of measurement. Slowly, doubts began to creep back into the statistics.¹ In the 1970s the psychologist Paul Meehl would worry that such brute approximations had stifled the development of the softer sciences and contributed to the mis-measurement of man. He would go on to elaborate twenty features of psychological science which made such measurement constructs inapt, unreliable and difficult to clearly falsify. This was progress!² He then showed that the methods of validating such constructs were the main culprit and the likely cause of psychology's implausible claims to concrete, replicable results. Around thirty years later some of the same replication issues would come to be called a crisis. Cournot was right to fear.

Seen from the perspective of a patient the difficulties are more urgent and stark. In 2019 Esme Yang would write with hope of the comfort given by a diagnosis, the knowledge that she was not "pioneering an inexplicable experience."³ To find yourself described in the DSM provides: a framework and the glimmer of a cure, a chance to slot yourself into a category, a community of care and a medicinal regime. Frustrating then when the categories themselves are in flux. Diagnostic criteria arrayed over page after page attempting to capture the elusive core of a psychological dysfunction. Descriptions are derived from clinical interviews correlated and meshed with genetic markers, then poorly mapped to a family of ailments. Neither clearly bipolar disorder nor manic depression, schizophrenia, psychosis or schizo-affective disorder. But an idiosyncratic presentation (as they all must be), uniquely felt and individually suffered.

¹ S.M. Stigler. *The Seven Pillars of Statistical Wisdom*. Harvard University Press, 2016

² Paul E Meehl. Theoretical risks and tabular asterisks: Sir Karl, Sir Ronald, and the slow progress of soft psychology. *Journal of Consulting and Clinical Psychology*, 46, 1978

³ All remarks on Schizophrenia in the following draw on the precise and personal essays in

Esme Weijun Wang. *The Collected Schizophrenias*. Penguin, 2019

Wang's collection of schizophrenias defy easy taxonomy "because there are just so many different ways in which people can develop a syndrome that looks like schizophrenia ... as we now define it." The task then becomes one of coping with uncertainty and adjusting expectations. It's not so exact a science that you can measure twice and cut once. The measurement schemes change and you need multiple cuts. You measure out the impact of treatments and the trade-offs - what's non-negotiable and how many side-effects are acceptable? What works for you versus what's advised by the professionals. Sterile cost-benefit considerations become suddenly dramatic and life changing.

...and Expected Value

There is an algorithm beloved by bureaucrats. An unsung hero of administrators and accountants. An algorithm both ubiquitous and under appreciated. It's pivotal for nearly every business and informs the actions of tech firms and policy makers the world over. It is only mildly hyperbolic to say that understanding this formula unlocks wealth and power. The algorithm lies at the heart of online A/B testing, all policy analysis, sound strategy and poker play.

$$EV(O)_p = p_1u(o_1) + p_2u(o_2) + \dots + p_ku(o_k)$$

The expected financial value of a random process is just the sum of the utility (typically dollar outcomes) weighted by their probabilities. Outcomes can vary from deals of cards, to customer transactions, election results, continued sanity. Pascal can argue sincerely that such considerations compel even belief in God. The infinite downsides of hell at any likelihood ought to compel even the cynical sceptic. But the formula, glossed as a rule for rational action, merits your attention for more mundane wagers too. If you intend to maximise your expected value, the meaning of probability is not an idle concern.⁴ While statistics are often tortured to rubber stamp decisions and probabilities are abused to fit policy prescriptions with false precision, the crisp clarity of the rule has an enduring appeal that promises to sift the murky swamps of Big Data. It's a scalpel that anyone can wield to parse the syntax of statistical jargon and carve answers from an abstract space of probabilities. "What's my expected return? My likely life-quality?" - a simple question, with a surprisingly complex answer.

Whenever expected utility is used as a metric of profit there is a risk of pivot, a point where an explanatory model of rational expected action is substituted for more obscure black-box optimisation tech-

⁴ For a typical example of how *EV* is used to express decisions under uncertainty see chapter 7 in the textbook [Barber, 2012], but the paradigm owes much of its influence to decision theoretic work of Leonard Savage in [Savage, 1954]

niques. The focus switches from modelling the consumer or the patient to modelling returns based on the consumer and the firm's expected value. Swapping a customer catering model for a customer-as-commodity perspective focused on profit in aggregate and customer acquisition. The tendency is common because it's easier and profit often overwhelms all other priorities, but the loss of understanding typically amounts to a longer term net-loss. Shareholders take comfort from increased profit and algorithms deployed at scale, but it's rare that any single algorithm actually or always maximises the available profit. This dynamic enforces a kind of inescapable see-saw motion where the consumer modelling exercise goes through a constant feedback loop. A good model (informal or formal) of human needs and wants feeds a better a predictive model of individual action. When the latter fails we go back to the utility curves and the algorithm of expected value because it is (if not reliably predictive) a rich and deeply explanatory model of human action under uncertainty.

Decision theory is an abstract formalism which purports to account for how you ought to reason under uncertainty, conditional on knowing your own mind, what you want and the likelihood of those outcomes. Wrestling with a diagnosis of schizophrenia you weigh your future in the face of sickness while knowing your mind could fail in the act, knowing, perhaps, you're not yourself. Schizophrenics are involuntarily detained if a medical professional deems them to have lost sufficient "insight"⁵, but the stated dysfunction assumes that these kinds of insight should be transparent when the psychiatric crisis is resolved. We'll see that the level of insight assumed by decision theory is, at best, hard won and far from transparent.

⁵ "The mind has been *taken over*. The mind has *lost the ability to make rational decisions*" in [Wang, 2019] pg58

Part I: Probability Measures

Probability: Two sides of a Coin

Probability emerged slowly and with a dual aspect. On one tradition probability refers to the long run tendency of a random process, on another probability is construed as the degree of belief in an outcome. On the first (frequentist) interpretation a probability distribution has certain fixed theoretical characteristics: as in a uniform probability distribution of a fair coin where all outcomes are equally likely, or as with the normal distribution where most outcomes cluster symmetrically about a central average. On the second (Bayesian) reading the characteristics of the probability distribution are learned from the data. The controversy centres around the fact that it's unclear how a frequentist could ascribe probabilities to unique events.

Without appeal to a large set of observations (or known theoretical distribution) the claim that an event appears frequently or infrequently is not well defined. Consequently, tabulations of probability appear inappropriate for claims of unique or rare events. In contrast the Bayesian is content to ascribe probabilities to any all partial beliefs no matter how specific. For the Bayesian, the probability calculus is a set of edicts about how to rationally manage and modulate your beliefs. So it's acceptable to have a probabilistic belief in rare cases so long as you update those probabilities with new data when available.

These two approaches are united by the Law of Large numbers which states that as the size of our sample increases our sample average will converge to the expected realisation of the theoretical process.

$$\frac{1}{N} \sum_{i=1}^N O_i \text{ converges to } E(O) \text{ as } N \text{ approaches } \infty$$

In this graph (Figure 1) we have fixed a Poisson distribution with a mean of 4.5 and can see three examples of how consecutive averaging from the increasing sample sizes results in a closer and closer convergence to the (true) population mean. Though common knowledge today, in 1650 "the very concept of averaging is [new]... and most people could not observe an average because they did not take averages."⁶ Systematically grappling with the implications of observations is a somewhat recent human endeavour - one which is far from perfected. This tendency is now fundamental to the interpretation of probability. Take a game with fixed and fair odds and we see that repeated play will converge over time because of characteristics which govern the process. Dice are a homely example. In the wild we never know the characteristics which cause the observed spread of outcomes, but such is the influence of gambling on probability, that we assume there exists a stable pattern to be gamed. Partially this is pragmatic, the maths is more tractable if we can assume a well behaved underlying process. The results are compelling: The Normal (Bell Curve) distribution, the Poisson distribution the Bernoulli distribution (to name a few) are all rightly famous. Their shapes are characteristics of innumerable random processes. The distributions cleanly circumscribe and corral likely patterns of events.

But the gambling paradigm clouds the fact that in practice we start on the left side of the law of large numbers (with samples) and we often start with small numbers resulting from a unknown number of data-generating processes. Well behaved probability distributions are rare beasts; a tiny fraction of the world's arbitrary menagerie. The

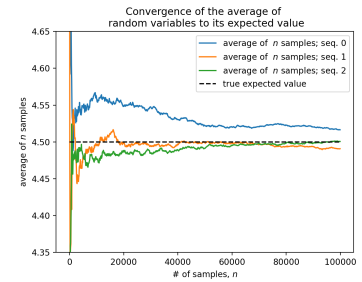


Figure 1: Convergence with large samples

⁶ Ian Hacking. *The Emergence of Probability: A Philosophical Study of Early Ideas about Probability, Induction and Statistical Inference*. Cambridge University Press, 2 edition, 2006

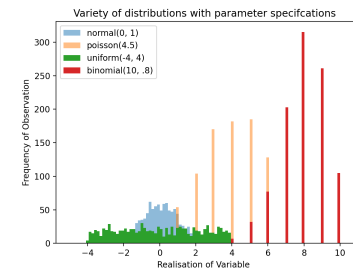


Figure 2: Some theoretical distributions with parameters

fundamental question in probability is not whether probability is a measure of belief or frequency - it is whether we can safely assume that the underlying process adheres to a known model? The Bayesian focus is to try and learn from the new data the expected characteristics of the underlying process, while the Frequentist tries to gauge the accuracy of their assumptions about the underlying process. Both are attempts to validate the structure of the model's theoretical distribution to inform inference. If we can't validate a model, we're better learning what we can from the sample, trusting to wide confidence intervals and worst scenario planning. But in all cases when we need to make a decision, the following questions are inescapable : What are your expectations based on? How do they figure in our choices, and can we use them to improve our outcomes?

Small Worlds and Statistical Inference

Anissa Weier and Morgan Geyser, later diagnosed with schizotypy, were in 2017 trailed for the attempted murder of their friend Payton. While believing themselves to be acting at the behest of Slender Man, they were "willing to forgo even friendship for the sake [their] version of unreality"⁷. Just as in the depth of an obsession reality can be warped by a psychotic focus, a statistical model will accentuate some parts and ignore others.

⁷ "The Slender Man, the Nothing and Me" in [Wang, 2019]

This narrowness can have devastating effects if deployed without care. But prediction is a visceral need, unavoidable, it precedes the sophistication of probability in any order of analysis. Without some regularity between X, Y we can only interpret their collisions as timorous noise. But even when there is a better than arbitrary correlation between X, Y we'll insist on ascribing a measurable probability to their association. Heuristics will kick in, and confidence will grow out of proportion to the evidence. So whether we view probabilities as a measure of credibility or frequency, the focus is always on a process which in reality reveals a pattern under repetition. Even tenuous connections can be enough to cause havoc. Models, more often than not, are deliberate simplifications of complexity, attempts to formulate a unobserved process within a mathematical structure. Lego-like versions of the world are built and rebuilt, in which we bash parts together to see what reliably sticks. Forget about the pieces we don't own. Count the pairs of blocks that: match, wedge, smash, click-into-place or break, then draw out the ratio of success and the spread of outcomes. Parse out the details of how reds go with greens, and blues with yellows. This is your sample distribution. The expected gain depends on both this

uncertainty, measured in this small world, and the finer points of statistical inference.⁸ . The danger of shrinking your world comes when you mistake the map for the territory and act on that delusion.

Small worlds are machines for figuring out expected values of a statistical process. We shrink the parameter set to better measure the variance and flux throughout the system. For any hypothetical system there can be multiple plausible approximations of the underlying process which need to be assessed comparatively on their “goodness” of fit. We iterate through new and improved versions, each an attempt to make a conjectural connection between X, Y mathematically precise. But once built, they embed the errors and assumptions made in their design. McElreath dubs them Golems, primed and then loosed on the world, insensitive to subtlety and context they perform only as instructed. Consider how a basic regression model tries to predict an outcome Y as linear function of some observed features X :

$$Y = \text{const} + \beta X + \epsilon$$

where ϵ is a random variable representing the error (or noise). A modest notational device for disaster. While const, β are parameters estimated by an optimisation process to ensure the equation fits the data as neatly as possible. In (Figure 3) below we have a series characterised by change. After the first shock we can refit the model so that the line tracks well with the evolving data. After the second shock we try another refit, but the range the and variance of the data makes our basic model a poor fit, i.e. the data no longer exhibits a linear relationship. This presents three examples of error in the modelling process: (i) it’s difficult to identify (in the moment) those changepoints in the data which reflect structural change, (ii) the linearity assumptions that go into the model are sound but the parameters need be re-estimated based on the new data and (iii) the third linear model is simply a terrible match for the pattern in the data.

Every model is a guess as to the implicit order in apparent noise. Sometimes there is no order, and other times the patterns is too subtle for a dumb model to capture. In practice you never really know whether a single new error stems from a misfit but appropriate model or an entirely inappropriate model. As we increase our number of sample fits we hope to better approximate the true linear process (if any) generating the data. Imagine now that the data points in Figure 3 are repeatedly re-speckled over the canvas. We can refit a new model for each set of scattered data points and each refit gives us a new sample values for const, β . If the underlying data generating process is stable, then the parameter fits will converge to the correct

⁸ A [small world] is... completely satisfactory, only if it is actually a microcosm, that is, only if it leads to a probability measure ... that can be written down explicitly pg 88

L. Savage. *The Foundations of Statistics*. Wiley, 1954

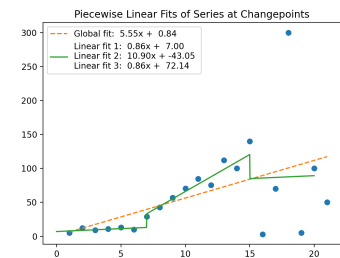


Figure 3: Three samples with starkly different parametrisations

values of $const, \beta$; correct in the sense that they can be used to draw the line of best fit for the data. A statistically stable process is one that can be modelled with errors ϵ normally distributed around 0, so that the model will be *correct on average* because $E(\epsilon) = 0$. Our predictions will overshoot in some cases but on the whole the errors up and down will cancel each other out.⁹ Forecasting with the parameters of best fit minimises our forecast errors because the fluctuations are stable about the centre of the line. These are the required assumptions for a process to exhibit the tendency of regression towards the mean. If they're not met, we will see poor parameter estimates and wild swings away from the linear path. The fundamental statistical assumption here is about the properties of our mistakes! The model is less plausible if our judgements are made in the grip of a delusion.

⁹ "Typically, the assumptions in a statistical model are quite hard to prove or disprove, and little effort is spent in that direction. The strength of empirical claims made on the basis of such modeling therefore does not derive from the solidity of those assumptions. Equally, these beliefs cannot be justified by the complexity of the calculations... These observations lead to uncomfortable questions"

David A. Freedman. *Statistical Models and Causal Inference: A Dialogue with the Social Sciences*. Cambridge University Press, 2009

Modeling: Improper Assumptions and Skewed Expectations

Below we build two sampling distributions based on different models of an underlying processes. One in which the errors are independent, normally distributed around 0 and in the other the errors are correlated in a sine-wave like pattern, increasing and decreasing periodically. This is akin to the difference between measuring error when predicting the heights of randomly sampled people versus predicting the sales volumes on randomly selected days of the week. A random sample of daily sales risks clumping weekends together and skewing the expected values. No such risk exists when sampling from independent individuals.

Build True Models

```
N = 100000
X = random.uniform(0, 20, N)
independent_err = random.normal(0, 10, N)
corr_err = random.uniform(0, 10) + sin(np.linspace(0, 10*pi, N)) +
sin(np.linspace(0, 5*pi, N))*2 + sin(np.linspace(1, 6*pi, N))*2

Y_corr = -2 + 3.5 * X + corr_err
Y = -2 + 3.5 * X + independent_err

population = pd.DataFrame({'X': X, 'Y': Y, 'Y_corr': Y_corr})
```

Sample from Data

and build smaller models

```
fits = DataFrame(columns=['iid_const', 'iid_beta', 'corr_const',
'corr_beta'])
```



```

for i in range(0, 10000):
    sample = population.sample(n=100,
                               replace=True)
    Y = sample['Y']; X = sample['X']
    Y_corr = sample['Y_corr']
    X = add_constant(X)
    iid_model = OLS(Y, X)
    results = iid_model.fit()
    corr_model = OLS(Y_corr, X)
    results_2 = corr_model.fit()
    row = [results.params[0], results.params[1],
           results_2.params[0], results_2.params[1]]
    fits.loc[len(fits)] = row

fits.boxplot()

```

In the case with independent errors the expected value for our parameter estimates match almost exactly the true values of the process. In the second model with correlated errors the parameter estimate for our constant is 4.9 which is significantly different from the true value of -2, and will lead to systematically skewed predictions. Statistical models are just algebraic equations where we use regular sampling to solve for Y from X . Because Y is also a random variable the regression model encodes a conditional expectation result.

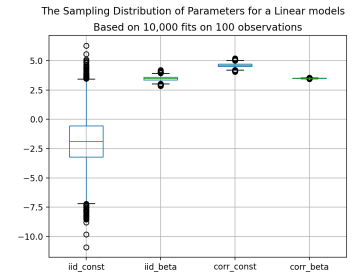


Figure 4: The expected realisations for $\beta, const$ with different errors structures



Figure 5: Error Distribution for the two models on a random sample of 1000 observations

$$E(Y_i | X_i = x)$$

For fixed values of X , the predictions Y_i can be spread in a pattern enforced by the various ways we can realise the linear function with estimates for β and $const$. But the regression model selects the best parameter values to minimise the squared prediction error and represent the conditional expected distribution of Y .¹⁰ The consequent point predictions for Y are always expected values, skewed by the

¹⁰ "The statement that regression approximates the [Conditional Expectation Function] lines up with our view of empirical work as an effort to describe the essential features of statistical relationships without necessarily trying to pin them down exactly" p38 -

Joshua D. Angrist and Jörn-Steffen Pischke. *Mostly Harmless Econometrics: An Empiricist's Companion*. Princeton University Press, 2008

how the parameters are realised from sample data as much as by poor choices in model design and predictive features. So too then any measures of expected utility based on these models or inferences from these distributions.

Frequentism: Inference from Expected Frequency

Making inference from a model is delicate thing. Even simple cases come with controversy. Count the number heads in a series of 5 successive coin flips, then repeat the process 1000 times and you'll arrive at a proportion which characterises that process. If it's a fair coin the long run expected result will be half the number of your coin flips. If the coin is weighted you might have as few as 0. This is the binomial distribution, and it really shines when you're trying to gauge fairness. If a process is biased, the distribution will be skewed. We can use this fact for inference. Consider a dispute over whether the game was rigged.

$$H_0 : \text{true proportion of heads} = 0.4$$

$$H_1 : \text{true proportion of heads} = 0.5$$

Take (H_0) as given then if we observe a sequence:

$$(3 \text{ in } 5) : H, H, H, T, T$$

what does it say about the possibility that we're being hustled? If the coin is biased, then the count of heads in repeated sampling will reflect a clear bias. For any new data we can check if the data is consistent with the data generated by the biased coin. The pattern of reasoning is straightforward (i) make some assumptions about the structure of the random process under investigation, (ii) tease out the consequences of these assumptions (iii) evaluate the incoming data against these consequences to see if you need to revise your assumptions. The frequentist asks, does the data looks weird given the assumed shape of our probability distribution?¹¹

In this instance the shape of the binomial distribution defined by a 0.4 biased coin allows for significantly greater than 5% chance for observing the above sequence. So we do not have enough reason to reject (H_0) at the traditional threshold. By design the assumed distribution builds in characteristics of long-run variance of the process, and the slim threshold for rejection is designed to minimise incorrect rejections of (H_0). We should remain suspicious that we're being conned. However, with a low number of observations the sample distribution is unlikely to be properly representative. This makes

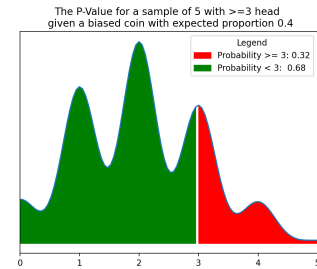


Figure 6: The Binomial Distribution

¹¹ " [I]n statistical terms H_0 [the null hypothesis] refers to a probability model and the very word 'model' implies idealization. With a very few possible exceptions it would be absurd to think that a mathematical model is an exact representation of a real system... We use the term to mean that in the current state of knowledge it is reasonable to proceed as if the hypothesis was true." pg 31

David R. Cox. *Principles of Statistical Inference*. Cambridge University Press, 2006

even small p-value thresholds unreliable. We cannot blindly take a sample poll to imply the spread or volatility of a population, and with low or un-representative samples it's hard to justify any kind of inference from expectation, since we are not in a position to justify the choice of the null model either! If your hypothesis is both derived and validated against small samples, you risk being swayed by recent observations. More fundamentally the notion of statistical significance usually cannot falsify the hypothesis under consideration. The sheer number of auxiliary variables that you might need to control for, makes the practical task of definitively rejecting the null almost impossible. There are too many imagined ways in which the auxiliary conditions, sufficiently modified, would have resulted in observations that corroborate the null model. This problem is especially acute in psychological science where the auxiliary contingencies of designing a measurement scale and checking diagnostic criteria are almost always questionable. So the null hypothesis is not confirmed, but not refuted either, it is just preserved in useless stasis.¹²

Bayesian Inference: Inference to Expected Value

If instead we use probability to calibrate our beliefs, then we can be more explicit in our assessment of (H_0) , (H_1) . Let's assume that our prior beliefs about whether the game is rigged is 50/50. Then we evaluate the two hypothesis using Bayes's rule for incorporating our prior belief and the data. The Bayesian asks whether our hypothesis is a good explanation of the data compared to alternatives. How, upon observing the data, should we view our hypothesis?

$$p(H_i|Data) = \frac{\overset{\text{prior}}{p(H_i)} \overset{\text{likelihood}}{p(Data|H_i)}}{\sum_{i=1}^K \underset{\text{evidence}}{p(Data|H_i)p(H_i)}}$$

where $1 \leq i \leq K$ spans the ways in which the data could have been realised across all competing hypotheses.¹³ Then, in our toy example, we have:

$$\frac{p(H_1|3in5)}{p(H_0|3in5)} = \frac{\frac{.5 \cdot .23}{.5 \cdot .32 + .5 \cdot .23}}{\frac{.5 \cdot .32}{.5 \cdot .32 + .5 \cdot .23}} = \frac{.57}{.42}$$

which would lead us to infer that the coin was fair. The really radical move in the Bayesian setting is that you're allowed to ascribe a probability to any event regardless of whether there is any long-run sequence to observe. You may know nothing about your opponent or the coin, but for Bayesians this is no bar to assigning suspicion in the form of expected probability, so long as you act in

¹² "[I]t did not get integrated into the total nomological network, nor did it get clearly liquidated as a nothing concept, it did not get killed or resurrected or transformed or solidified; it just kind of dried up and blew away..." pg807 in [Meehl, 1978]

¹³ "When a piece of evidence E is produced in a court investigating the guilt G or innocence I of the defendant, it is not enough merely to consider the probability of E assuming G; one must also contemplate the probability of E supposing I. In fact, the relevant quantity is the ratio of the two probabilities. Generally if evidence is produced to support some thesis, one must also consider the reasonableness of the evidence were the thesis false. Whenever courses of action are contemplated, it is not the merits or demerits of any course that matter, but only the comparison of these qualities with those of other courses." in

D.V. Lindley. The analysis of experimental data: The appreciation of tea and wine. *Teaching Statistics*, 1993

accordance with the axioms of probability and weigh the probabilities accordingly. In particular it promotes the direct comparison of competing hypotheses conditional on the evidence. It's this free choice of prior which can seem arbitrary and unmotivated or even paradoxical, but in practice probabilities are rarely ascribed entirely without reason and it's frankly irresponsible to ignore those reasons.

Neither the Bayesian or Frequentist analysis ends with these simple calculations, both should continue to probe the limits of each hypothesis. We'd have to consider things like sample size, sensitivity testing, model performance, the cost of errors and appropriateness of the priors. The point is just that there are reasons for dispute. Bayesian inference acts like a logic engine for evidence, whereas the frequentist approach is more focused on diagnosing the possibility of error. In general they are complementary methods, and when they conflict the assumptions should be scrutinised. The frequentist evaluation of our biased coin is very sensitive to the choice of hypotheses, while the Bayesian approach is influenced by the choice of prior. Why set up a significance test against assumed cheating rather than assumed fairness? Why attribute equal weight to both hypotheses? Why use a 5% threshold if you're concerned about systematic cheating? This example shows the heart of the conflict in the dual aspect of probability. There is enough latitude in the manner in which we set up a probability model that the mathematics of inference can yield inconsistent results. Both offer strategies for managing uncertainty, but both approaches come with baggage and in practice not all tests are equally taxing. Consider a more concrete example in the Bayesian spirit.

An Example: Expected Website Returns

Websites and apps collect traffic and log interactions. Your details are captured and pulled into vast aggregates of consumer data. I can route and re-route your trajectory across an online environment. Applying the same pressures to tens of thousands of others, we can trace out how the topology of particular sites throw up speed bumps on the customer's journey. Imagine we're running a website which aims to funnel customers through to a number of different purchase plans. The historic patterns are relatively stable with only 10% of customers dropping out of our conversion funnel on a daily basis. We can sample actions online (Figure 7) under differing pressures with a view to evaluating expected values of repeated coercive prompts.

Assume the particular values for each plan, then the expected value of customer journey is just: $p_1\$ (o_1) + p_2\$ (o_2) + p_3\$ (o_3) + p_4\$ (o_4) =$

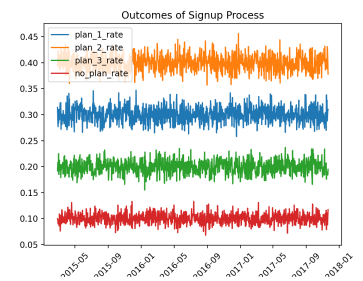


Figure 7: Stable long run Sign Ups

$.3 * \$10 + .4 * \$7 + .2 * \$12 + .1 * \$0 = \$8.20$. Now imagine there was a change to the website and we observe the following pattern (Figure 8) for the next 20 days. The change was made on the hypothesis ($H+$) that it would bring a positive boost to revenue. How much more positive? A slight expected increase makes it harder to conclusively reject the $H+$ even in the fact of contrary indicators.

What is the new expected value? Have we decisively falsified $H+$? From the frequentist point of view the macro distributional properties haven't significantly changed. But given what we know about the change to the website it would be foolish to accept such a static distributional assumption. Looking only at the small sample of new data, the variance will be large and the estimates of rates of sign-up for each plan will be unstable. Following the Bayesian paradigm we can condition our expectations on the new data, the old data or all the data. The below graph illustrates the spread in values expected revenue calculated on different slices of our data using Bayes Rule. Using a large number of observations, the influence of our priors are minimal and washed out by the data, giving us a strong point estimate with low variance stable around 8.2, but since the recent data involves a step change, we might be better off ignoring the old data. But we can also see that if we condition our expectations only on the new data with different priors drawn from the past data or hope, we can positively bias our expectations. Suppose we're naive and accept either the optimistic prior or retain a frequentist approach, and accept that the website change is associated with a slight drop in revenue, do we revert to the old website or try to explain the dip by contingencies of the market and preserve our test for another thirty days?

This pattern is not rare. Nearly all substantial decisions are made with small samples in circumstances where past behaviour is not a guide. Past behavioural patterns are exactly what we're trying to avoid or change. If you want to know whether the change on your website will drive a material change in financial revenue, you won't have long run patterns to rely on, and it's an open question on how to weight the new data. If you want to judge the long term consequences of a new symptom the same limitation of information applies. All models smuggle-in a host of statistical assumptions and these can range from reasonable to absurd. Even when reasonable they're only supported by large sample sizes, and most questions of interest are driven by novelty (or specificity) that short circuits appeal to robust patterns of history. Reasoning from small samples is common, best done with caution and plenty of caveats, but better reason than not. Expectations should be modified accordingly.

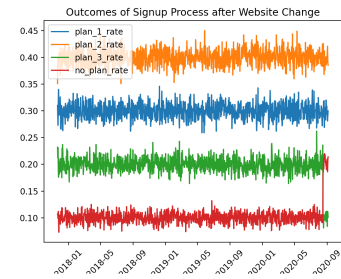


Figure 8: Abrupt increase in dropouts

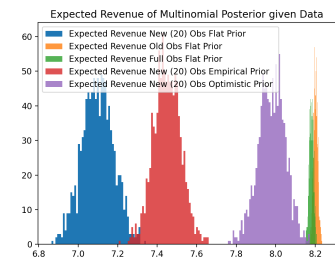


Figure 9: Expected revenue differs by choice of prior and data

Part II: Utility Curves

The Stakes: From Utility to Indifference

Our views of probability can flex up and down in response to facts, but it's less clear how our estimates of utility change. Too much of a good thing often tends to the bad. So we dabble, sample and share. In pursuit of variety we swap our goods, shunning stale options in favour of the novel exchange. For a given good we can differ in our appetites but it's relatively straightforward to find the point where - one more donut is one too many. While it can be a bit unclear how we should measure utility, once we've decided on a metric the mathematical characteristics are meaningful. If the scale is donuts, we can infer aspects of your attitudes from your acquisition and enthusiasm for donuts. In most cases we're interested not just in your pursuit of pastries, but how you'd be willing to trade for those pastries.

We seek competitive advantage for our own produce to balance the cost owed to the skills of others. This coordinated compromise lies at the core of maximising subjective utility in a market, but at the limit some scenarios do not admit any admixture of goods. Not all babies can be cut in half. But in most cases a consumer will try to optimise their bundle of goods over an entire marketplace, preserving enough on one key good; money, to remain liquid. So, to a first approximation our utility estimate would seem to be a multivariate function.

$$u(\mathbf{g}) = f(g_0, g_1 \dots g_n)$$

There are number of ways we can specify a utility function as seen in Figure 10, but a typical example is the Cobb-Douglas function.

$$u(\mathbf{g}) = g_0^{\alpha_0} g_1^{\alpha_1} \dots g_n^{\alpha_n} \text{ where } \forall i \sum \alpha_i = 1$$

Then taking the case of two goods g_1, g_2 we can determine an indifference curve where you would be willing to exchange quantities of g_1 for an agreeable amount of g_2 . The task is to express the value of a given good as priced in terms of the other goods. Set

$$\begin{aligned} u(\mathbf{g}) &= k = g_1^{\frac{1}{2}} g_2^{\frac{1}{2}} = (g_1 g_2)^{\frac{1}{2}} = \sqrt{g_1 g_2} \\ \Rightarrow k^2 &= g_1 g_2 \Rightarrow \frac{k^2}{g_2} = g_1 \end{aligned}$$

Using this formula we can express how the quantities of fair exchange vary based on a fixed utility value. This is not to say that these curves represent an actual or objectively fair price, just that

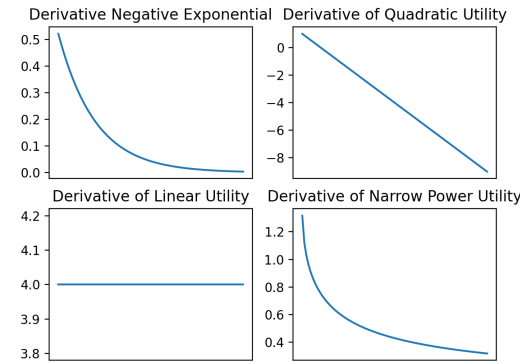


Figure 10: The Rates of Change of personal Utility

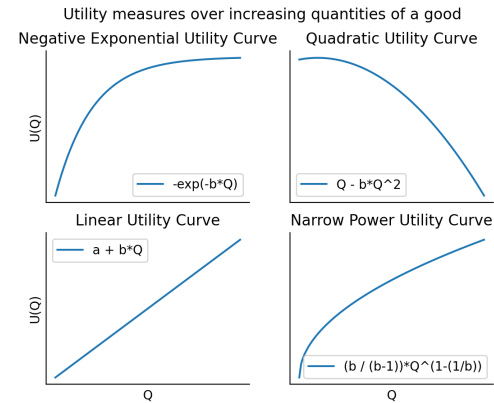


Figure 11: Consumer attitudes with differently satisfied appetites for a good

when measured in terms of our utility these are mappings of quantities of goods we would be happy to exchange. Your view of a fair price is encoded in your utility theory. It's at this point when utility theory can be said to verge on empirical science. If we can model your preferences as a utility function characteristic of some general attitude toward acquisition, we might also hope to be able to predict future trades and cater for individual desires.

Contrast the case of a schizophrenic's patient's utility as they try to weight aspects of their health. Even if you could measure each dimension simultaneously, what scale captures the worth of autonomy and measures the trade off against clarity, delusion and paranoia? How do you plan for children in the fear of what your genetics could seed? Do you fear the harm done through inheritance or the looming lifelong responsibility of care for an affected child? Even adoption foists onto the child the burden of having you as a parent. The disease calibrates the cost, but the equations are hard to solve.¹⁴

Optimising Utility

A further complication arises when we try to account for a consumer's budget or resources. The shape of the Cobb-Douglas function in Figure 12 shows that the utility surface is constantly increasing with our rate of acquisition. So without any constraints the consumer would not achieve satisfaction, but continue like a glutton consuming forever without cease. Add a budgetary constraint and natural trade-offs between desire and cost mean that we need to find the maximum point at which an indifference curve intersects with our budgetary line. Instead of solving the equation:

$$\text{Find } g_1, g_2 \text{ such that } u(g_1, g_2) = \lambda$$

we need to solve a constrained optimisation problem:

$$\text{maximize } u(g_1, g_2) \text{ subject to } \text{cost}(g_1, g_2) = \lambda$$

This style of problem can be approached with the method of Lagrange multipliers. If we let:

$$L = g_1^{\frac{1}{2}} g_2^{\frac{1}{2}} - \lambda(2g_1 + 3g_2 - 40)$$

where 2 and 3 are the unit cost of the respective goods, and 40 is our total budget. This λ is our Lagrangian multiplier - a term used to re-express the algebra of our equation as a function of the consumer's capacity to spend.

We can discover where utility is maximised when the gradient of the "curve" can be set to zero. This is the theory behind the "hill

Cobb Douglas Utility Curve for two goods

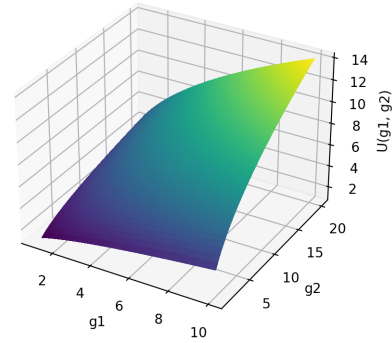


Figure 12: A consumers utility curve for combinations of two goods

¹⁴ See "The Choice of Children" in [Wang, 2019]

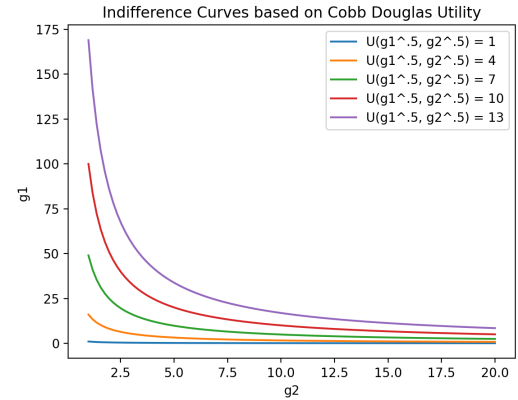


Figure 13: A range of indifference curves without budget constraints.

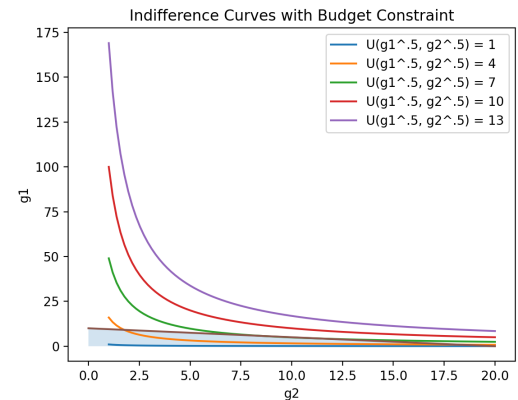


Figure 14: A range of indifference curves with budget constraints.

climbing" algorithms of gradient descent. When the curvature of the "slope" has plateaued i.e. is zero, then we've reached a maximum or minimum in the multivariate space of the function. As before we want to use this fact to express the implicit function of g_1 in terms of g_2 , but this time including the constraints on our budget.

Box .1: Lagrangian Multiplier

$$\nabla L = dL/d\mathbf{g} = \left(\frac{\partial u(\mathbf{g})}{\partial g_1}, \frac{\partial u(\mathbf{g})}{\partial g_2} \right) = \left(\frac{1}{2}g_1^{-\frac{1}{2}}g_2^{\frac{1}{2}} - 2\lambda, \frac{1}{2}g_2^{-\frac{1}{2}}g_1^{\frac{1}{2}} - 3\lambda \right) = \mathbf{0}$$

$$\Rightarrow \lambda = \frac{1}{4}g_1^{-\frac{1}{2}}\sqrt{g_2} = \frac{4}{25}g_2^{-\frac{1}{2}}\sqrt{g_1}$$

$$\Rightarrow \left(\frac{1}{4}\right)^2 \frac{1}{g_1} g_2 = \left(\frac{4}{25}\right)^2 \frac{1}{g_2} g_1 \Rightarrow \left(\frac{1}{4}\right)^2 g_2 = \left(\frac{4}{25}\right)^2 \frac{1}{g_2} g_1^2 \Rightarrow \left(\frac{1}{4}\right)^2 g_2^2 = \left(\frac{4}{25}\right)^2 g_1^2$$

$$\Rightarrow g_2 = \frac{16}{25}g_1$$

The same pattern holds for cases with more than two goods. We can express the value of given good g_n in terms of a function $f(g_1, \dots, g_{n-1})$. Then substituting this value into our constraint we get:

$$2g_1 + 3\left(\frac{16}{25}\right)g_1 = 40 \Rightarrow 2g_1 + 1.92g_1 = 40 \Rightarrow 3.92g_1 = 40$$

Proving the optimal settings are $g_1^* = 10.20$ and $g_2^* = 6.52$ and $\lambda^* = 0.20$

The above proof shows how we triangulate a consumer's view of any good as expressed through the medium of their utility function. But the method of Lagrangian multipliers is more than a mere algebraic trick. We can interpret the λ term as the rate of change of the consumer's utility as a function of the cost to our resources. How taxing is our treatment, how exhausted is your wallet? The proof is a little more involved, but the significance of this interpretation should be obvious. If we knew our consumers adhered to a particular style of utility function we could model how "price-changes" would impact their returns to utility and select strategies for maximum gain. The challenge lies in deriving a customer's utility profile.

Part III: Representation Theorems

Rational Preference: From Indifference to Utility

The core idea is that an agent's utility metric ought to reflect their preferences, so if we can elicit preference statements from our consumer, we should be able to construct their utility curve! One method suggests itself; first canvas a customer for their preferences or ob-

serve their preference as expressed by purchases. Then map the maximal and least preferred options to convenient polarities. For instance:

$$g_1 \succ g_2 \succ g_3 \succ g_4 \succ g_5$$

where:

$$u(g_1) = 0 \text{ and } u(g_5) = 1$$

then each of the intermediary options can be measured in the interval between 0 and 1. However, not all relations map to preference structures e.g. there are an infinite number of simple ordinal mappings that would work, but a strict ordering does nothing to convey the degree of feeling associated with each option. Extra constraints need to be placed on the utility metric if it is to reflect the properties of a genuine preference relation. Our preferences need to respect certain axioms of rationality.

Box .2: Individual Preference Structures

Let R denote a binary relation over a set of states S . We can place a variety of conditions on the R preference relation:

- **Reflexivity** $\forall \phi \in S : \phi R \phi$
- **Completeness** $\forall \phi, \psi \in S : (\phi \neq \psi) \rightarrow (\phi R \psi \vee \psi R \phi)$
- **Transitivity** $\forall \phi, \psi, \chi \in S : \phi R \psi \wedge \psi R \chi \rightarrow \phi R \chi$
- **Anti-Symmetry** $\forall \phi, \psi \in S : (\phi R \psi \wedge \psi R \phi) \rightarrow \phi = \psi$
- **Asymmetry** $\forall \phi, \psi \in S : \phi R \psi \rightarrow \neg(\psi R \phi)$
- **Symmetry** $\forall \phi, \psi \in S : \phi R \psi \rightarrow \psi R \phi$
- **Acyclic** $\forall \phi_1 \dots \phi_j (\phi_1 R \phi_2 \dots R \phi_j \rightarrow \neg(\phi_j R \phi_1))$

Definition: (Individual Preference Structure) $\langle S, \succeq, \succ, \sim \rangle$ is a (relational) structure where \succeq is a weak preference ordering on the set S if \succeq is a transitive, reflexive and complete. While \succ is a strict preference ordering if $\forall \phi, \psi : \phi \succ \psi \Leftrightarrow \phi \succeq \psi \wedge \neg(\psi \succeq \phi)$. Finally, we have \sim as an indifference relation if $\forall \phi, \psi : \phi \sim \psi \Leftrightarrow \phi \succeq \psi \wedge \psi \succeq \phi$

Definition: (Choice Functions) We let C be a choice function when $C : S \mapsto S^* \subseteq S$ just when $S^* \neq \emptyset$ unless $S = \emptyset$. For example we have a choice function with respect to a preference structure when $C^\succeq(S) = \{\Phi \in S : \forall \Psi \in S, \Phi \succeq \Psi\}$

The technique pursued by Von Neumann and Morgenstern is to calibrate utility scales based on decisions made about offered bets. Each individual good can be assessed against a simple win-loss lottery between the two most extreme outcomes. If the consumer is indifferent

between the sure prospect of the good and a fixed odds lottery on their most (and least) preferred outcomes, they've implicitly weighed their utility of the good.

$$\forall g_i \exists p : g_i \sim [p \cdot g_1, (1-p) \cdot g_5] \rightarrow u(g_i) = p$$

So whenever we are indifferent between a sure thing and a win-loss lottery over the best and worst outcomes we have implicitly chosen the utility of the good on a 0-1 scale. In this manner we can construct a utility curve across the entire range of options to reflect an underlying preference relation.

Von Neumann-Morgenstern Representation: Decision Under Risk

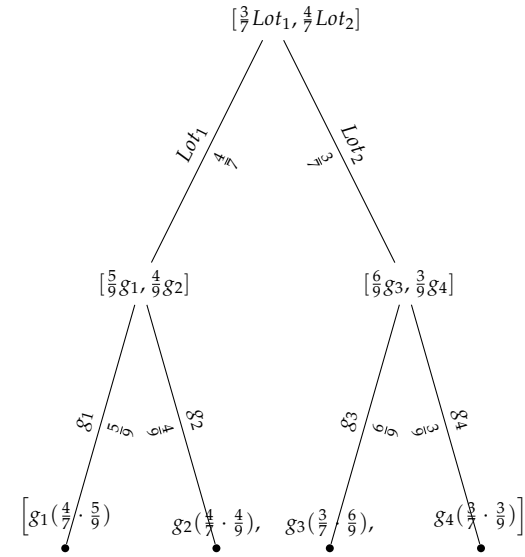
The most famous result in decision theory is von Neumann and Morgenstern's Representation theorem. It shows, using the technique discussed above, how expressed preferences (which adhere to certain axioms of rationality) can track with a utility measure. As such the agents can be interpreted as making choices to maximise their expected utility. But the theorem is limited to decisions over well-defined lotteries, and as such makes a poor model for general choice under uncertainty where the odds are approximate, unknown or unclear. Nevertheless the theorem serves as an alternative explanation for the tendency to make decisions based on expected value.

Theorem 1 (vNM's Representation Theorem) *If an individual i's preference relation \succeq is transitive, complete and satisfies:*

1. (Continuity): $\forall g_1, g_2, g_3 : (g_1 \succeq g_2 \succeq g_3) \rightarrow \exists v \in [0, 1] \wedge g_2 \sim_i [vg_1, (1-v)g_3]_{Lot}$
2. (Monotonicity): *If $v_1, v_2 \in [0, 1]$ and $g_1 \succ g_2$ then $([v_1g_1, (1-v_1)g_2]_{Lot} \succeq [v_2g_1, (1-v_2)g_2]_{Lot}) \Leftrightarrow v_1 \geq v_2$*
3. (Reduction of Compound lotteries): *Each compound lottery $[q_1Lot_{p_1}, \dots, q_nLot_{p_n}]$ reduces to a simple lottery where each good $(1, \dots, k)$ is weighted across all branches of the nested decision tree $[(q_1p_1^k + q_2p_2^k \dots + q_np_n^k)g_k, \dots, (q_1p_1^{k-1}, \dots)g_{k-1} + \dots (q_1p_1^1, \dots)g_1]_{Lot}$ by the usual rules of conditional probability for branching events such that $\widehat{Lot} \sim Lot$*
4. (Independence) *If $\widehat{Lot} = [q_1Lot_1, \dots, q_jLot_j \dots q_nLot_n]$ and $L_j \sim M$, then $\widehat{Lot} \sim \widehat{Lot}' = [q_1Lot_1, \dots, q_jM \dots q_nLot_n]$*

then $\exists u_p : \widehat{Lot} \mapsto Val$ where $u_p(\widehat{Lot}) = p_1u(g_1) + \dots + p_ku(g_k)$ and $u(\widehat{Lot}) \geq u(\widehat{Lot}^) \Leftrightarrow \widehat{Lot} \succeq \widehat{Lot}^*$ so that u represents \succeq unique up to*

Figure 15: Reduction of Compound lotteries as probability trees



a positive linear transformation.

For a well defined and fixed probability function p over the goods g_1, \dots, g_n the above axioms of rationality are sufficient to define a sensible utility function based on an agent's expressed preferences.¹⁵ The thought gives hope to the idea that you would be able to predict an individual's actions in any environment where you knew both their preferences and the objective probabilities at play. This is the basic model for understanding poker play - the probabilities are generally known and it just remains to determine the game theoretical dynamics, assuming the other players act consistently and intelligently to pursue rational preferences.

¹⁵ "The point is that there is no need to assume or philosophize about, the existence of an underlying subjective utility function, for we are not attempting to account for the preferences or the rules of consistency. We only wish to devise a convenient way to *represent* them". p32

R.D. Luce and H. Raiffa. *Games and Decisions: Introduction and Critical Survey*. Dover Publications, 1989

Bolker-Jeffrey Representation: Decision Under Uncertainty

There is an altogether different view of what we're doing when we pursue expected utility. Again, we may try to elicit an agent's utility function from their preferences, but in addition we can set up the axioms of rationality so as to derive a probability function based on the expressed desires. This is a Bayesian approach to what's going on when we act to maximise expected utility - one which emphasises the dynamic and subjective nature of the attributed probabilities.¹⁶ In particular this means we can attribute probabilities and preferences across boolean combinations of arbitrarily complex propositions instead of lotteries and compound lotteries. The propositions express an individual's belief and the probabilities can be thought of as the strength of the belief.

¹⁶ Richard C. Jeffrey. *The Logic of Decision*. University of Chicago Press, 1983

For this interpretation to work we need some extra structure to represent simple algorithmic rules for composition of belief. Taking simple cases and aggregating or combining them in a way which adheres to obviously sensible procedures in the base case, and generalises across a total range of arbitrary complexity.

Box .3: Boolean Algebras

Definition (Boolean Algebra) is a relational structure $\Omega = \langle S, \wedge, \vee, \neg, \top, \perp \rangle$ such that the following axioms hold:

- **Associativity** $\alpha \wedge (\beta \wedge \gamma) = (\alpha \wedge \beta) \wedge \gamma$ and $\alpha \vee (\beta \vee \gamma) = (\alpha \vee \beta) \vee \gamma$
- **Commutativity** $\alpha \vee \beta = \beta \vee \alpha$ and $\alpha \wedge \beta = \beta \wedge \alpha$
- **Absorption** $\alpha \vee (\alpha \wedge \beta) = \alpha$ and $\alpha \wedge (\alpha \vee \beta) = \alpha$
- **Idempotence** $\alpha \wedge \top = \alpha$ and $\alpha \vee \perp = \alpha$
- **Normality** $\alpha \vee \neg \alpha = \top$ and $\alpha \wedge \neg \alpha = \perp$
- **Distributivity** $\alpha \wedge (\beta \vee \gamma) = (\alpha \wedge \beta) \vee (\alpha \wedge \gamma)$ and $\alpha \vee (\beta \wedge \gamma) = (\alpha \vee \beta) \wedge (\alpha \vee \gamma)$

Example: The classical propositional calculus forms a boolean algebra where the elements $\top, \perp \in S$ are interpreted as Truth and Falsity. The set of well-formed formulas of propositional logic are: $\alpha \mid \alpha \wedge \beta \mid \alpha \vee \beta \mid \neg \alpha \mid$ and we have an interpretation function $\llbracket \cdot \rrbracket : \Omega \mapsto \{1, 0\}$ across the signature of the language in the usual truth functional manner.

$$\llbracket \alpha \wedge \beta \rrbracket = \llbracket \alpha \rrbracket \wedge \llbracket \beta \rrbracket$$

$$\llbracket \alpha \vee \beta \rrbracket = \llbracket \alpha \rrbracket \vee \llbracket \beta \rrbracket$$

$$\llbracket \neg \alpha \rrbracket = \neg \llbracket \alpha \rrbracket$$

The mapping defines an implication relation $\alpha \models \beta \Leftrightarrow \llbracket \alpha \vee \beta \rrbracket = \llbracket \beta \rrbracket \Leftrightarrow \llbracket \beta \wedge \alpha \rrbracket = \llbracket \alpha \rrbracket$
 $\Leftrightarrow \llbracket \alpha \rrbracket \leq \llbracket \beta \rrbracket$ defines an ordering.

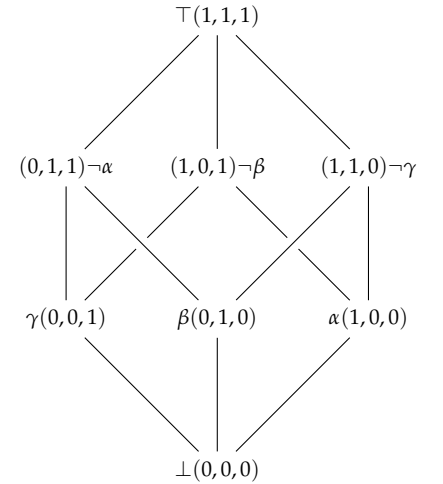
From Neutrality to Desire

One of the issues with eliciting a utility curve with appeals to bets over lotteries stems from the stigma associated with gambling. An alternative approach, more in the spirit of Bayesian philosophy is to try to elicit the desirability of a prospect by situating it between two polarities and repeatedly seeking a third prospect, midpoint between the two, which is half as desirable by construction. The method originally proposed by Frank Ramsey relies on the idea that we express preferences over a boolean algebra of propositions and we can gauge utility by appeal to an "Ethically Neutral" proposition *Neutral* - one which if it exists is such that for all other prospects α we're utterly indifferent between:

$$(Neutral \wedge \alpha) \sim \alpha \sim (\neg Neutral \wedge \alpha)$$

. We can gauge desire by offers of repeatedly refined contracts based on an ethically neutral proposition. This sequence of offers can be used to construct a utility curve. This is Bayesian in spirit because it allows for the expression of a probability measure for any proposition even if there is no canonical probability distribution over the considered outcomes.

Figure 16: Boolean Algebra of Propositions



First observe that

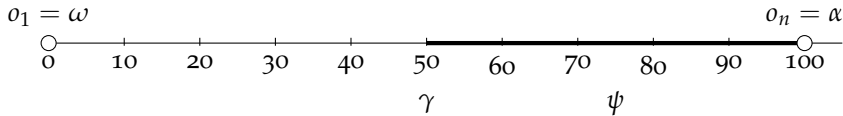
$$\begin{aligned}
 & [\alpha \text{ if } \textit{Neutral}, \omega \text{ if } \neg \textit{Neutral}]_{\textit{contract}_1} \\
 & \sim [\omega \text{ if } \textit{Neutral}, \alpha \text{ if } \neg \textit{Neutral}]_{\textit{contract}_2} \\
 & \Rightarrow u(\textit{contract}_1) = u(\textit{contract}_2) \\
 & \Rightarrow EU(\textit{contract}_1) \\
 & = u(\alpha)p(\textit{Neutral}) + u(\omega)(1 - p(\textit{Neutral})) \\
 & = u(\omega)p(\textit{Neutral}) + u(\alpha)(1 - p(\textit{Neutral})) \\
 & = EU(\textit{contract}_2) \\
 & \Leftrightarrow p(\textit{Neutral}) = 0.5
 \end{aligned} \tag{1}$$

Then we can take any two extremes $u(\alpha) = 1, u(\omega) = 0$ and we can use our test for indifference to situate any third proposition γ on a desirability scale since:

$$\begin{aligned}
 & [\gamma \text{ if } \textit{Neutral}, \gamma \text{ if } \neg \textit{Neutral}]_{\textit{contract}_1} \\
 & \sim [\alpha \text{ if } \textit{Neutral}, \omega \text{ if } \neg \textit{Neutral}]_{\textit{contract}_2} \\
 & \Leftrightarrow EU(\textit{contract}_2) = u(\alpha)\frac{1}{2} + u(\omega)\frac{1}{2} = .5 \\
 & = u(\gamma) = EU(\textit{contract}_1)
 \end{aligned} \tag{2}$$

Repeating this step we can find a contract on a sure-thing ψ for which we're indifferent between:

$$\begin{aligned}
 & [\psi \text{ if } \textit{Neutral}, \psi \text{ if } \neg \textit{Neutral}]_{\textit{contract}_1} \\
 & \sim [\alpha \text{ if } \textit{Neutral}, \gamma \text{ if } \neg \textit{Neutral}]_{\textit{contract}_2} \\
 & \Rightarrow u(\psi) = .75 \\
 & \dots \textit{etc}
 \end{aligned} \tag{3}$$



Repeating this process indefinitely we can refine our utility scale as exactly as we please, by repeatedly finding prospects with a utility precisely on the mid-point between two poles. If these measures adhere to certain basic constraints of rationality regarding consistency of utility we can show how a Bayesian agent can be seen to maximise their expected utility when making decisions under uncertainty. But unlike the Von Neumann representation theorem, for a Bayesian the probability function over prospects is not unique. We can have multiple pairs $\langle p, u \rangle$ which are representative of an individual's preference ordering \succeq without converging on the particular probabilities ascribed by one individual. This is precisely the content of the following theorem.

Theorem 2 (Bolker Representation Theorem) Let $\mathbb{B} = \langle \Omega, \succeq \rangle$ be Bolker structure if Ω is an atomless Boolean algebra and \models forms an implication relation over Ω , while \succeq is complete, transitive, continuous over $\Omega \setminus \perp$ and the following hold:

1. (Impartiality) Suppose $\alpha \sim \beta$ and $\exists \gamma (\neg(\gamma \sim \alpha))$ such that $\alpha \wedge \gamma = \perp = \beta \wedge \gamma$ and $\alpha \vee \gamma \sim \beta \vee \gamma$. Then $\forall \gamma (\alpha \vee \gamma \sim \beta \vee \gamma)$
2. (Averaging) If $\alpha \wedge \beta = \perp$ then $\alpha \succeq \beta \Leftrightarrow \alpha \succeq \alpha \vee \beta \succeq \beta$

Then there is a probability measure and utility (desirability) metric $\langle p, u \rangle$ on Ω such that if the following axioms hold:

- (A0) $p(\top) = 1$
- (A1) $p(\alpha) \geq 0$
- (A2) $\alpha \wedge \beta = \perp \rightarrow p(\alpha \vee \beta) = p(\alpha) + p(\beta)$
- (A3) $u(\top) = 0$
- (A4) $\alpha \wedge \beta = \perp \wedge p(\alpha \vee \beta) \neq 0$ implies

$$u(\alpha \vee \beta) = \frac{u(\alpha)p(\alpha) + u(\beta)p(\beta)}{p(\alpha \vee \beta)}$$

it follows that

$$u(\alpha) \geq u(\beta) \Leftrightarrow \alpha \succeq \beta$$

and there is another such set of functions $\langle p^*, u^* \rangle$ if and only if u^* is a fractional linear transformation of u i.e. $\exists a > 0$ and $\exists c, cu(\alpha) > -1$

$$p^*(\alpha) = p(\alpha) \cdot (cu(\alpha) + 1)$$

$$u^*(\alpha) = \frac{au(\alpha)}{cu(\alpha) + 1}$$

The mathematical machinery used to prove this result is a little more involved, ranging over every possible boolean combination of beliefs measured on three axes: preference, probability and desirability. In addition to the usual probability axioms, (A3) and (A4) tie subjective probability and subjective utility together. The axiom (A3) works to normalise the utility scale so that no sure prospect has any positive utility. This, in a sense, enshrines the requirement that there is only a utility to novel information. While (A4) ensures that the utility of any disjunction is the weighted average of the ways in which it can occur. More importantly it implies:

$$\begin{aligned} u(\alpha \vee \neg\alpha) &= u(\top) = p(\alpha)u(\alpha) + p(\neg\alpha)u(\neg\alpha) \\ &= p(\alpha)u(\alpha) + u(\neg\alpha) - p(\alpha)u(\neg\alpha) \end{aligned}$$

$$\begin{aligned} \Rightarrow u(\top) - u(\neg\alpha) &= p(\alpha)u(\alpha) - p(\alpha)u(\neg\alpha) \\ \Rightarrow p(\alpha) &= \frac{u(\top) - u(\neg\alpha)}{u(\alpha) - u(\neg\alpha)} \text{ if } u(\alpha) \neq u(\neg\alpha) \end{aligned}$$

Which confirms how the relationship between probability of a given proposition can be expressed in terms of the desirability or utility of the same proposition and it's negation. This is a view of probability profoundly different from measure of risk we ascribe to players calculating pot-odds in poker. It is not a fixed unique distribution determined by observation across repeated sampling, and consequently much harder to model. The probabilities reflect the dynamics and eccentricities of individual beliefs and the agent is seen as maximising their subjective expectations. Given how the average consumer cannot then be modelled with respect to a fixed reference probability distribution, you might despair of ever predicting an individual's actions. Fortunately crowding promotes conformity and what seems mysterious at the micro level becomes clearer at the macro scale.

Part IV: Machine Learning and the Individual

Customer Representation: A Segmentation Approach

The most influential approach to automating the representation theorems above stem from the work of Daniel McFadden. His analysis of the BART transport system in San Francisco and the algorithmic approach he took to estimating the preferences of the San Francisco residents. He used this analysis to accurately predict the uptake in the rail-users within the city thereby showing that a sound understanding of the incentives and pressures on city infrastructure can influence the preferences of the citizenry. The core insight treats utility as a latent factor driving demand in the market. The utility is some function of product and consumer's properties, perhaps mostly driven by price

$$utility = \mathbf{X}'\beta + v$$

and market share is an expression of that utility

$$demand_A = utility_A > 0$$

So we can use the observed facts of market demand to estimate β coefficients which determine subjective utility ranking across the market. The revealed preference assumption says that we can predict the purchase if the utility of the good is positive.

$$Pr(demand_A = 1) = utility > 0$$

$$\begin{aligned}
&= Pr(\mathbf{X}'\beta + v > 0) \\
&= Pr(v > -\mathbf{X}'\beta) \\
&= 1 - F(\mathbf{X}'\beta)
\end{aligned}$$

where F is the distribution of the unobserved random variable v . The challenge is using the correct distribution as this feeds the method of statistical estimation of the parameters β . So in the case of a choice over two goods, we have two equations and the utility of the product A is estimated as its positive difference over the reference product B.

$$U_{i,A} > U_{i,B}$$

just when

$$\mathbf{X}'_{i,A}\beta + v_{i,A} > \mathbf{X}'_{i,B}\beta + v_{i,B}$$

or

$$v_{i,A} - v_{i,B} > -(\mathbf{X}'_{i,A} - \mathbf{X}'_{i,B})\beta$$

but then the probability of demand is just

$$\begin{aligned}
&Pr(demand_A = 1 | \mathbf{X}'_{i,A}, A, \mathbf{X}'_{i,B}, B) \\
&= Pr\left(v_{i,A} - v_{i,B} > -(\mathbf{X}'_{i,A} - \mathbf{X}'_{i,B})\beta\right) \\
&= Pr\left(-v_{i,A} - v_{i,B} < (\mathbf{X}'_{i,A} - \mathbf{X}'_{i,B})\beta\right) \\
&= F\left((\mathbf{X}'_{i,A} - \mathbf{X}'_{i,B})\beta\right)
\end{aligned}$$

Which in practice tend to be estimated with logistic regression model in the binary case and the multinomial logistic regression in model when there are more goods to be considered. We take one of the product classes as a reference class evaluate the relative probability of demand.

$$Pr(y_A = 1 | \mathbf{X}_{i,A}, \dots, \mathbf{X}_{i,j}) = \frac{\exp(\mathbf{X}'_i \beta)_{i,A}}{1 + \sum_{j=1}^{j=N} \exp(\mathbf{X}'_i \beta)_{i,j}}$$

and we loop through each the products defining their respective worth in terms of the reference product. This is a strong restriction called **The Irrelevance of Independent Alternatives**, as it bakes in the notion that our preferences are consistent and transitive. If we prefer C to B and B to T then we ought to prefer C to T too. The benefit of the assumption is that it allows us to infer a utility ranking

metric by computing all the pairwise alternatives to a given product. The multinomial logit distribution is a convenient measure for discrete choice problems because it allows us to express our preference for each product on a 0-1 scale. McFadden's work used this technique to predict accurately the uptake in rail travel on the BART system in San Francisco, by first fitting a multinomial model on locations with a rail system to derive the coefficient weights of influence for a range of demographic factors around house size, home-ownership, income and location on demand for rail.

Dep. Variable:	Choice	No. Observations:	11624
Model:	MNLogit	Df Residuals:	11614
Method:	MLE	Df Model:	8
Date:	Mon, 01 Mar 2021	Pseudo R-squ.:	0.1071
Time:	22:58:06	Log-Likelihood:	-4499.7
converged:	True	LL-Null:	-5039.6
Covariance Type:	nonrobust	LLR p-value:	9.100e-228

Choice=car	coef	std err	z	P> z	[0.025	0.975]
home	1.0460	0.128	8.159	0.000	0.795	1.297
HHSIZE	0.2411	0.046	5.234	0.000	0.151	0.331
income	0.2235	0.021	10.829	0.000	0.183	0.264
urban1	1.7447	0.150	11.619	0.000	1.450	2.039
density	-0.0847	0.006	-13.677	0.000	-0.097	-0.073

Choice=rail	coef	std err	z	P> z	[0.025	0.975]
home	0.3757	0.140	2.689	0.007	0.102	0.650
HHSIZE	-0.0495	0.051	-0.973	0.331	-0.149	0.050
income	0.1658	0.022	7.466	0.000	0.122	0.209
urban1	-0.0949	0.167	-0.569	0.569	-0.422	0.232
density	0.0202	0.007	3.107	0.002	0.007	0.033

These modelling efforts give us a sense of how the relative utility of a given good is seen across the market by a range of consumers, but there are complications. In particular, most products come with a price which is correlated with the error term in linear equation driving the utility model. This confounds the clean estimation of the parameters. Combined with the questionable notion that the linear preference order represents an actual consumer preference, makes the modelling effort somewhat unreliable as a predictive tool.

PCA and Segmentation

Yet, treating the customer as commodity with a limited range of behaviour puts us back in the business of sampling. We treat each interaction as instance of behaviour fluctuating around a predictable pattern. This is comforting because it's familiar and seemingly absolves us about attributions of utility and the burden of interpretation. This is illusory. While there is a wide range of segmentation methods which can be applied to the task of classifying both customers and products. These classification schemas are vital inputs for any recommendation algorithm. They cluster individuals based on a wide array of features, which is to say that they simplify the question of expected action. Instead of asking how might Rebecca, (aged between 18-25, from Spain, with a history of frugal purchases) react to a new sales promotion, you can ask about the conversion rate of the young female demographic. Depending on the task and the nature of the clustering algorithm, you might end up bucketing Rebecca with Sven (overweight male, history of lavish spending from Sweden) if, for example, their historic email open rates were similar. The responsibility for vetting each clustering schema lies with the user of the algorithm, but usually knowledge of the problem is enough to put some kind of context on the structural patterns unearthed by the algorithm.

Customer Features				
Customer ID	custDesc0	custDesc1	...	CustDescN
1	3.2	4.5	...	10
2	5.2	4.3	...	8.2
3	5.6	4.2	8.5
4	7.5	4.6	...	12

Table 1: *

Ideally we would parse out our data into the most relevant descriptive features, but in lieu of domain knowledge we can apply some data compression techniques such as principle components analysis to extract latent features in the data. These techniques are dangerous when applied without domain knowledge or some kind of supervision. A technique, in the same vein as PCA, called factor analysis tries to construct latent factors from correlations in the observed features. Historically this was abused to measure "intelligence" as a latent factor driving performance on aptitude tests. So while we should be wary of over-inflating artefacts of the data, the technique is useful for finding structure. Starting from the covariance matrix of the scaled customer data :

$$\text{Cov}(\mathbf{X}) = \text{cov}[X_i, X_j] = E[(X_i - E[X_i])(X_j - E[X_j])]$$

```

X = df_customer[[x for x in df_purchases.columns
if 'customer_desc' in x]]
X_std = StandardScaler().fit_transform(X)
## Covariance Decomposition
cov_mat = np.cov(X_std.T)
eig_vals, eig_vecs = np.linalg.eig(cov_mat)
## Explained Variance
tot = sum(eig_vals)
var_exp = [(i / tot)*100 for i in sorted(eig_vals, reverse=True)]
cum_var_exp = np.cumsum(var_exp)

```

For any observed set of customer attributes, their covariance matrix can be factorized (or decomposed) into a matrix product of the eigenvectors and eigenvalues. These components can in turn be analysed to express the manner in which each of the original features contributes to the new construct. Each component captures a proportion of the total variance observed in the original covariance matrix: If we're lucky, a small number of the eigenvectors (principle components) can be shown to explain the variance in the data (as in Figure 17) and thereby represent a complex customer problem in a reduced dimensional space. Each customer is visualised as a complex of observed characteristics along the axes of the principle components in, for example, a two dimensional plane. Furthermore, we can overlay an interpretation on the components by associating group identifiers to portions of the plane. In the dataset pictured in Figure 19 we've tried to represent the data in three groups inferred by a k-means clustering algorithm over the original observed features. This is not the most appropriate categorisation as can be seen by the manner in which 9 clear cohorts are speckled across the space. In other words there is more diversity in our customer base than the clustering is capable of expressing. That's not to say all diversity needs to be captured. If our purpose is make the broadly correct action (e.g. sell, promote, no-action), the niceties can often be ignored. A Silhouette score analysis shows that one of our classes is too crude, collecting dissimilar customers together. Nevertheless, the expected value of the customer is then something like the probability of purchase (however approximated) for the given customer multiplied the the net revenue of the average purchase in the same group. It's these kind of idealisations and conceptual economies, which seemingly reasonable at the time, lead to systematic algorithmic skew. Unrefined constructs encourage brutish bulk actions across the arrayed field of customers.

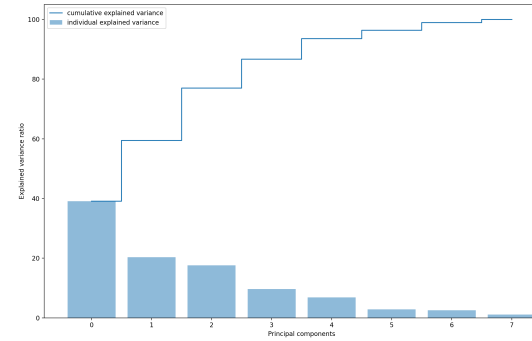


Figure 17: Principle Component Analysis

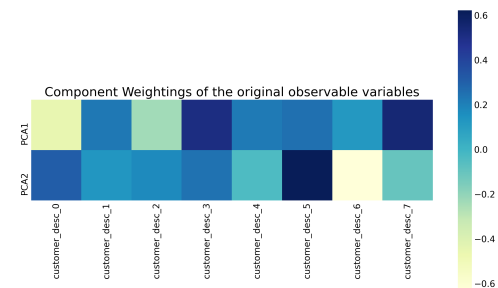


Figure 18: Principle Component Weightings by Observed Customer Features

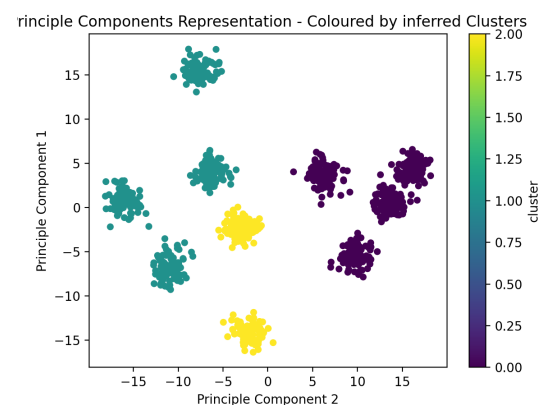


Figure 19: Principle Component Representation of Customer Clusters

$$EU^{seg} = \sum_{i=1}^{i=n} p^{seg_i} u^{seg_i}$$

Note the level of abstraction here! In an effort to understand the customer as a predictable entity we've gone from a set of concrete but complex customer descriptions through a matrix decomposition representation of their variant behaviours, then focused on a weighted sum construct of the observed behaviours that in some sense best "explain" the diversity of the behaviours in the original data. Concluding with a group-categorisation of how that behavioural construct can be interpreted or acted upon. The grouping is key. If the complexities of the model is a wound demanding attention. The segmentation is the interpretable overlay which is sold to your management, it's the salving balm of targeted advertisement sold to investors.

So in practice there is some constraints of "plausibility" applied to these kinds of models. But once the model is live in productions it becomes hard to correct for anything other than degrading performance on the metric that matters most; expected value.

Conclusion: Construct Criticism

The theory of expected value has been around for a long time. It has been criticised and corrected, adjusted and refined. We've seen two broad species of justification for the rule: (1) a random process will converge around a stable mean in the long-run by the law of large numbers, so better to maximise that mean, (2) the procedure is justified by appeal to the representation theorem and the various axiom schemes of rationality. In practice, both justifications flounder. We can rarely justify the time or resources required for long run convergence since most questions are more urgent, and there are regular examples of humans violating the principle of expectation maximisation when faced with a decision. A number of apparent paradoxes where following the paradigm would result in reputedly irrational results. Moves and counter-moves in the debate. Straw-men are built and burnt, but the back and forth is beyond the scope of this blog post.¹⁷ Nevertheless even with radically different approaches to modelling customer preference, concerns about expected value will permeate all proposed solutions. It is a metric, explanation and a soothing, bloodless formula.

Yet the model constructs which feed the formula are not neutral; on the one hand you may question the sample data, on the other there are fair questions about the design choices that go into construct-

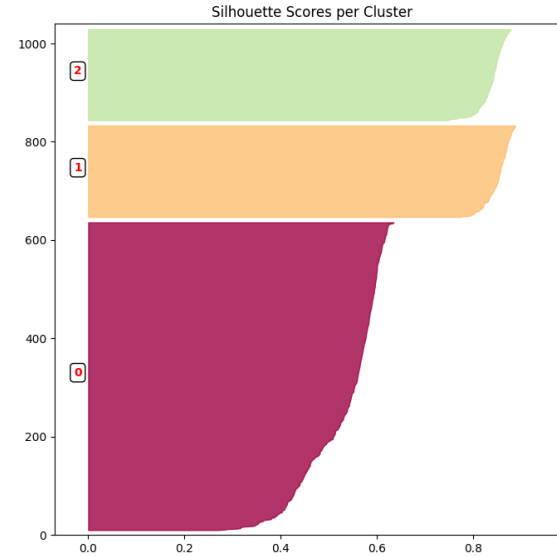


Figure 20: Cluster Validation by Silhouette scores - a customer measure of similarity within and across the available clusters.

¹⁷ A good discussion can be found in Ken Binmore. *Rational Decisions*. Princeton, 2011

ing the models. Did your data capture the correct features? Were you able to extract the important structural regularities? Did we run the experiment for long enough? Was the sample biased? Is the attributed probability still valid, are the utility curves steep enough? Does the clustering scheme make sense? When the model design is filtered through the expected value metric, the subtleties of those choices get glossed over and obscured. The creative leaps and structural assumptions calcify, become heuristics and accrue advocates and devotees. Analysts masquerade as oracles claiming as wisdom formulas carved in clay on crumbling architecture. The irony here is sharp. In both the Von Neumann and Jeffrey's formulation, decision theory is an inescapably constructive project - directly responsive to the polled preferences of dynamic individuals across a variety of different domains, subject to differing axioms as appropriate. Yet, stasis is the inevitable result from cleaving too tightly to a simple metric like expected value. Improvements are made and measured in minor swings around the sink-hole of local optima, if you're lucky. More likely, in a market context, the simple minded pursuit of a single KPI leaves you vulnerable to exploitation. Without regular attention these model constructs lose relevance and compound the errors of the past. The simple rule is often too simple.

References

- Joshua D. Angrist and Jörn-Steffen Pischke. *Mostly Harmless Econometrics: An Empiricist's Companion*. Princeton University Press, 2008.
- David Barber. *Bayesian Reasoning and Machine Learning*. Cambridge University Press, 2012.
- Ken Binmore. *Rational Decisions*. Princeton, 2011.
- David R. Cox. *Principles of Statistical Inference*. Cambridge University Press, 2006.
- David A. Freedman. *Statistical Models and Causal Inference: A Dialogue with the Social Sciences*. Cambridge University Press, 2009.
- Ian Hacking. *The Emergence of Probability: A Philosophical Study of Early Ideas about Probability, Induction and Statistical Inference*. Cambridge University Press, 2 edition, 2006.
- Joseph Y. Halpern. *Reasoning about uncertainty*. MIT Press, 2005.
- Richard C. Jeffrey. *The Logic of Decision*. University of Chicago Press, 1983.

D.V. Lindley. The analysis of experimental data: The appreciation of tea and wine. *Teaching Statistics*, 1993.

R.D. Luce and H. Raiffa. *Games and Decisions: Introduction and Critical Survey*. Dover Publications, 1989.

Paul E Meehl. Theoretical risks and tabular asterisks: Sir Karl, sir Ronald, and the slow progress of soft psychology. *Journal of Consulting and Clinical Psychology*, 46, 1978.

L. Savage. *The Foundations of Statistics*. Wiley, 1954.

S.M. Stigler. *The Seven Pillars of Statistical Wisdom*. Harvard University Press, 2016.

Esme Weijun Wang. *The Collected Schizophrenias*. Penguin, 2019.