

Ignorance and Indifference: Decision-Making in the Lab and in the Market*

Sergiy Verstyuk[†]

This version: 22 January 2018

Abstract

Economic agents face an evolving, non-ergodic environment. The corresponding permanently undersampled “population” distribution naturally permits unseen, rare events. The principle of indifference, implemented via the methods of parameterization invariance and maximum entropy, provides a disciplined, rational approach to learning, inference and decision in this underinformed setting. Canonical economic problems of choice under uncertainty from both the micro (binary lotteries) and macro (consumption-investment) domains can be formulated in terms of dynamic online learning when new gambles and new regimes are regularly encountered. The implied Bayesian updating under invariant ignorance priors allows to reverse-engineer and rationalize, in a mutually consistent way, the Allais paradox/prospect theory’s probability distortions identified in laboratory experiments as well as the equity premium/risk-free rate, non-monotone pricing kernel and portfolio underdiversification puzzles observed in financial markets.

*Acknowledgements: Pietro Veronesi, George M. Constantinides, Doron Ravid, Matt Taddy; as well as Jörn Boehnke, Christian Julliard, Leonid Kogan, Nicholas Polson, Xiao Qiao, Philip J. Reny, Lawrence D. W. Schmidt, Harald Uhlig.

[†]Contact address: verstyuk@cmsa.fas.harvard.edu.

Contents

- 1 Introduction** **1**

- 2 Empirics** **3**
 - 2.1 Calibrations 3
 - 2.2 Experiments 9
 - 2.3 Additional empirical results 12

- 3 Theory** **13**
 - 3.1 Conceptual approach 13
 - 3.2 Lotteries in laboratory experiments 14
 - 3.3 Assets in financial markets 19

- 4 Discussion** **26**

- 5 Conclusion** **27**

1 Introduction

A classical problem in statistics and decision theory is accounting for the unseen events. Pierre-Simon Laplace (1812) posed it as a question about the chances the sun will rise tomorrow given it always had for the past 5000 years, and formulated his rule of succession as a candidate solution to this sunrise problem. Alan Turing and I.J. Good (1953) encountered this issue while decoding the German Enigma machine cypher, later exemplifying it with the problem of estimating the number of yet undiscovered animal species/vocabulary words, and relating it to computational linguistics more generally. Recent studies tackle this problem more rigorously, dealing with optimal coverage-adjustment and smoothing of the distributions, identifying their shapes and various statistics such as entropy (see Nemenman et al., 2002; Orlitsky et al., 2002; Paninski, 2003; Hausser and Strimmer, 2009; Valiant and Valiant, 2015). The somewhat older works of Keynes (1921), Jeffreys (1939, 1946) and, very prominently, Jaynes (1957a, 1957b) were concerned with fundamentally related problems.

The need to conduct inference about, estimate and account for unobserved classes of contingencies, the very fact of whose existence is unknown, also arises in economics. Consumers, investors, job-seekers do not live in a static world but face an evolving stochastic environment that must be learned and responded to in real time. In elementary cases, they regularly encounter new tasks of choice under risk, often being forced to make a decision after only cursory experience (e.g., consider binary lotteries with limited information about the success probabilities). In more complex cases, economic regimes change, and what agents had learned about the previous regimes is of little or no use (e.g., parameters driving the macro environment are perturbed and re-initialized). That is, in a non-ergodic world, economic agents are permanently struggling with “small samples” (at least in the sense of effective sample size).

The above non-ergodic, substantially undersampled environment entails permanent parameter and model uncertainty. The principle of indifference, broadly understood, offers ways of dealing with such ignorance in a disciplined manner: namely, the method of maximum entropy and the method of invariance to transformations/reparameterization. Within the Bayesian framework and for the basic univariate cases we consider they reduce to Jeffreys’ approach to assigning prior probability distributions (exploiting conjugacy, they also admit transparent analytical solutions).¹

We apply this indifference-motivated approach to the canonical problems of decision-making under uncertainty: primitive binary lotteries as well as the consumption-investment choice. When these problems are built upon (online) learning of the unknown parameters of Bernoulli and Normal probability distributions, Bayesian updating with invariant ignorance priors produces the posterior distributions of, respectively, Beta and Student’s t

¹See Jaynes (2003) for an extensive discussion; Kass and Wasserman (1996) provide an excellent review.

(mixture) forms that are biased towards quasi-uniformity. Moreover, in a (permanently) undersampled environment, the effect of such non-informative priors never washes away, amplifying the relatively low-density boundary/tail parts of the corresponding posteriors.²

On the empirical side, these theoretical results are consistent with observations on the behavior of human subjects in laboratory experiments as well as with the asset pricing and portfolio holding regularities in financial markets. In the case of binary lotteries, comparing the available sample and inferred posterior, or “true”, distributions produces what resembles the Allais (1953) “paradox” and the probability weighting transformations stipulated by the (cumulative) prospect theory of Kahneman and Tversky (1979, 1992); which suggests rational optimality-motivated foundations for (some aspects of) the purely descriptive prospect theory.³ In the case of the consumption-investment choice, the recognition that the observed sample distribution might differ from and not reflect in full the risks of the unobserved posterior (or “true”/“population”) distribution makes it possible to reproduce the equity premium and risk-free rate “puzzles” found in the post-World-War-II United States data (Mehra and Prescott, 1985; Weil, 1989) without the need to resort to implausible levels of risk aversion.

Moreover, not only are the above theoretical results in microeconomic/decision theory and in macroeconomics/finance obtained using fundamentally the same approach, but the corresponding empirical results are also mutually consistent: the incorporation of the probability weighting transformation measurements from the former allows to match the latter. This establishes the connection that contributes to a unified understanding of additional empirical regularities such as the non-monotone pricing kernel (or the implied volatility “smile”) and the portfolio underdiversification “puzzles”.

Thus, under the imperative of parameterization invariance and/or maximum entropy, the documented “biases” (i.e., the probability distortions and large risk-premia identified in the laboratory experiments and in the financial markets) are just manifestations of the participants’ perceptions about the challenging stochastic environment around them. In other words, people embrace uncertainty and try to make rational decisions in the face of uncertainty.

Lastly, the distinction and—in practice, permanent—divergence between the sampled data and the inferred posterior distributions reconciles the allegedly mutually inconsistent normative theoretical desiderata of standard von Neumann-Morgenstern axioms on the one hand (which, in our framework, apply to the latter) with the positive empirical observations of decision-making under risk on the other hand (which correspond to the former), thus killing two birds with one stone.

In addition to literature references included throughout the text, the studies closely

²By definition, non-ergodicity means that the sample statistics are not representative of the underlying process. So, learning from history only sharpens, but never dominates, the priors.

³A very selective list of literature measuring such probability weights consists of the older studies by Tversky and Kahneman (1992), Camerer and Ho (1994), Wu and Gonzalez (1996), Abdellaoui (2000) as well as the newer ones by Bruhin et al. (2010), Fehr-Duda and Epper (2012).

related to the present one are as follows. In microeconomics and decision theory, it is Steiner and Stewart (2016) with their (neural) noise-in-processing-information story (as well as the focusing on the costs of self-control Fudenberg and Levine, 2006; 2011). In macro-finance, primarily it is the inventive work by Weitzman (2007) with a Bayesian view on the “peso problem” (as well as Tsai and Wachter, 2016; and Veronesi, 2004); additionally see other studies on the role of investors’ beliefs, especially the relative-entropy based approaches (Hansen, 2014, contains a fresh overview). Separately, there are also careful empirical papers by Polkovnichenko (2005) and Polkovnichenko and Zhao (2013) on the interface of the two subfields. Our main innovation is an attempt to take a unifying perspective, as well as providing an optimality argument for the existing empirical findings.

2 Empirics

2.1 Calibrations

We start with an inspection of the standard consumption and portfolio choice model applied to U.S. consumption as well as government bill and stock returns data. The model is due to Lucas (1978), a general equilibrium macro-finance workhorse based on an endowment economy.

The top panel of Table 1 provides key macro-financial facts on the United States after World War II. These are the sample statistics observed by econometricians, they characterize the realized joint probability distribution via the moments of its marginals $h_R(R_{t+1})$, $h_C(C_{t+1})$ and $h_D(D_{t+1})$.⁴

In the bottom panel of Table 1, we run the model under different assumptions about (i) the unrealized and unobserved population distribution, parameterized by its proximity to the sample distribution in terms of mutual information distance κ (which is a convenient measure of the distributions’ similarity and uncertainty that does not rely on the existence of statistical moments; see MacKay, 2003, for a good introduction), as well as (ii) the investors’ attitude towards risk, parameterized by the coefficient of relative risk aversion γ . For the relevant probability density $\pi(\cdot)$, the risk-free rate there is calculated as

$$\mathbb{E}^\pi[R_0] := \frac{1}{\mathbb{E}^\pi[M_{t+1}]} \quad (1)$$

and the risk premium as

$$\mathbb{E}^\pi[R - R_0] := -\frac{\text{Cov}^\pi[R_{t+1}, M_{t+1}]}{\mathbb{E}^\pi[M_{t+1}]}, \quad (2)$$

where

$$M_{t+1} := \beta \frac{u'(C_{t+1})}{u'(C_t)} = \beta \frac{u'(D_{t+1})}{u'(D_t)} \quad (3)$$

⁴We abstract away from labor income, although empirically it constitutes a large share of the national income.

Table 1: Calibration Results

		$E^h[R_0]$	$E^h[R - R_0]$	$E^h[\Delta c]$	$E^h[\Delta d]$											
		$\sqrt{V^h[R_0]}$	$\sqrt{V^h[R]}$	$\sqrt{V^h[\Delta c]}$	$\sqrt{V^h[\Delta d]}$											
Empirical, h		0.88	8.23	1.94	2.43											
		1.26	17.05	0.99	4.23											
		1			2			3			4			5		
γ :	κ	$E^*[R_0]$	$E^*[R - R_0]$	$E^*[\Delta d]$	$E^*[R_0]$	$E^*[R - R_0]$	$E^*[\Delta d]$	$E^*[R_0]$	$E^*[R - R_0]$	$E^*[\Delta d]$	$E^*[R_0]$	$E^*[R - R_0]$	$E^*[\Delta d]$	$E^*[R_0]$	$E^*[R - R_0]$	$E^*[\Delta d]$
		$\sqrt{V^*[R_0]}$	$\sqrt{V^*[R]}$	$\sqrt{V^*[\Delta d]}$	$\sqrt{V^*[R_0]}$	$\sqrt{V^*[R]}$	$\sqrt{V^*[\Delta d]}$	$\sqrt{V^*[R_0]}$	$\sqrt{V^*[R]}$	$\sqrt{V^*[\Delta d]}$	$\sqrt{V^*[R_0]}$	$\sqrt{V^*[R]}$	$\sqrt{V^*[\Delta d]}$	$\sqrt{V^*[R_0]}$	$\sqrt{V^*[R]}$	$\sqrt{V^*[\Delta d]}$
Model, h	n/a	6.18	2.85	8.07	6.17	2.86	3.60	6.17	2.86	2.30	6.17	2.86	1.69	6.19	2.86	1.34
		n/a	17.05	17.01	n/a	17.05	8.41	n/a	17.05	5.59	n/a	17.05	4.19	n/a	17.05	3.35
Model, g	0.1	-6.32	16.25	8.07	-6.35	16.27	1.99	-6.36	16.28	0.89	-6.36	16.28	0.50	-6.35	16.28	0.33
		n/a	40.37	40.26	n/a	40.36	19.22	n/a	40.36	13.10	n/a	40.36	9.81	n/a	40.36	7.85
Model, g	0.2	0.36	8.78	8.07	0.34	8.79	2.87	0.33	8.80	1.66	0.33	8.80	1.16	0.35	8.80	0.89
		n/a	29.80	29.73	n/a	29.14	14.64	n/a	29.80	9.73	n/a	29.80	7.29	n/a	29.80	5.83
Model, g	0.3	2.65	6.38	8.07	2.64	6.39	3.16	2.64	6.39	1.92	2.64	6.39	1.37	2.65	6.39	1.07
		n/a	25.44	25.37	n/a	25.44	12.52	n/a	25.44	8.32	n/a	25.44	6.23	n/a	25.44	4.98
Model, g	0.4	3.80	5.21	8.07	3.79	5.22	3.30	3.79	5.22	2.05	3.79	5.22	1.48	3.80	5.22	1.16
		n/a	23.01	22.95	n/a	23.01	11.34	n/a	23.01	7.53	n/a	23.01	5.64	n/a	23.01	4.51
Model, g	0.5	4.48	4.53	8.07	4.47	4.54	3.39	4.46	4.54	2.12	4.46	4.54	1.54	4.48	4.54	1.21
		n/a	21.47	21.41	n/a	21.47	10.58	n/a	21.47	7.03	n/a	21.47	5.27	n/a	21.47	4.21

Notes: The top panel “Empirical” presents statistical moments for the real risk-free and market (excess) returns, the real per capita consumption growth as well as the real dividend growth. The data sample is U.S. 1948:Q1–2014:Q4, at quarterly frequency (see Appendix §C for a description of data sources). The bottom panel with “Model” rows presents the same statistical moments as above (the real per capita consumption growth is omitted being identically equal to the real dividend growth) for different combinations of the mutual information parameter and the coefficient of relative risk-aversion. The CRRA utility, $\beta = 0.99$, and the probability density of returns (h or g) are the models’ only inputs (see text for a detailed description). The measurement units are nats for the mutual information distance, and percentage points converted into annualized terms for the economic variables.

is a (model-induced) stochastic discount factor (SDF), also called a pricing kernel or the marginal rate of substitution between consumption at time t and time $t + 1$. Equations (1)–(3) let us solve for the levels of consumption and dividend growth implied by the above-mentioned distributional and risk-aversion assumptions.

As a baseline, we start with the empirically observed sample $h_r(\cdot)$; see the rows corresponding to “Model, h ” of Table 1. Specifically, we evaluate our model using as inputs the probability density of observed returns ($R \sim \log \mathcal{N}(\hat{\mu}_r, \hat{\Sigma}_r)$ for the values of $\hat{\mu}_r$ and $\hat{\Sigma}_r$ estimated in sample) as well as the discount factor $\beta := 0.99$ for the plausible values of the RRA coefficient γ (implemented at quarterly frequency). However, the computed solution is not consistent with the rest of the observed data. Most striking are the low values obtained for the risk premium $E^h[R - R_0]$ and the high values for the risk-free rate R_0 for the realistic choices of parameter γ . These are the manifestations of the classic “equity premium puzzle” (Mehra and Prescott, 1985) and “risk-free rate puzzle” (Weil, 1989). From such a result, we can infer that the probability distribution (and/or SDF specification) used in the actual investment decision-making differs from what we fed into the model.

Next, taking the distributions’ mutual information as a free parameter, we search for the value of κ that produces Σ_r that lets the model fit the empirical observations picked as testing criteria. In other words, we adjust the variance upward, trying to back-out the true Σ_r and thus the putative distribution $g_r(r)$. (Also, to ensure that $E^g[R_{t+1}] = E^h[R_{t+1}]$, log-Normality of returns requires us to adjust μ_r , so that in the end $\mu_r < \hat{\mu}_r$.)⁵

The results obtained in the same manner as those for “Model, h ” earlier, but using as a probability density input the adjusted density $g_r(\cdot)$, are placed in rows “Model, g ” of Table 1. (For comparison, the sample-based “Model, h ” is effectively a special case of “Model, g ” with $\kappa \geq \kappa^* := 3.4$.) In the vicinity of $\kappa = 0.2$ and $\gamma = 3$, the model results compare well with the data. The risk premium is slightly above 8% and the risk-free rate is slightly below 1% are close to empirically observed values. Of course, some words of skepticism are in order: the level of the risky asset’s return was used as the model input. However, the ultimate split within this given level into the layer of the risk-free return and a risk premium on top is dictated by the model. Moreover, the mean and standard deviation of consumption growth (as well as dividend growth, since in our model featuring just one risky asset they coincide) are at low and high single-digit values, respectively, which looks plausible in comparison to empirical yardsticks; and this time no direct relation to model inputs can beget any skepticism.⁶ The result of the adjustment

⁵Given that the true distribution $g_r(\cdot)$ is not known and what we work with is just its calibration/estimate, one might find it more apt to denote objects like Σ_r with double-hats, i.e. $\hat{\hat{\Sigma}}_r$.

⁶The true distribution’s mean consumption (or dividend) growth obtained in our calibration is 1.66%, which is not too far from empirical values of 1.94% (2.43%), adding plausibility. The corresponding standard deviation is 9.73%, which is naturally larger than 0.99% (4.23%) observed in the sample, but is comparable to alternative results. Specifically, it is lower than the 17% standard deviation of welfare equivalent consumption growth in Weitzman (2007); it falls in between the 4.24% and 14.84% standard

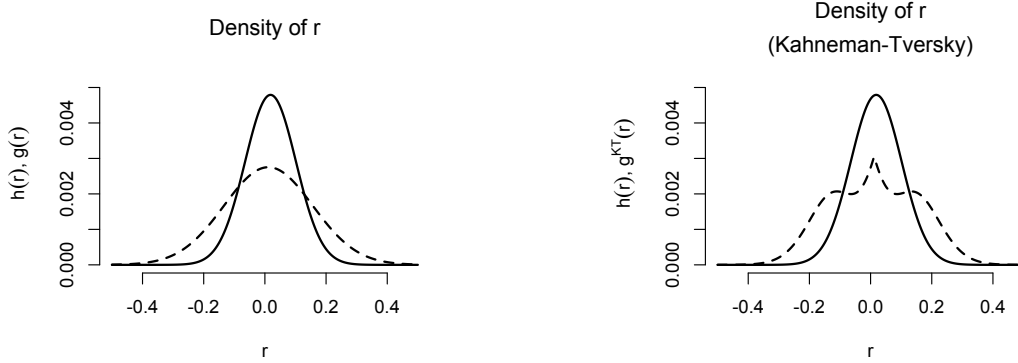


Figure 1: Empirical (solid) vs. adjusted (dashed) probability densities (parameterizations used are $\mathcal{N}(\hat{\mu}_r, \hat{\Sigma}_r)$ with $\hat{\mu}_r := 0.0184$ and $\hat{\Sigma}_r := 0.0069$ vs. $\mathcal{N}(\mu_r, \Sigma_r)$ with $\mu_r := 0.0114$ and $\Sigma_r := 0.0210$, corresponding to “Model, g ” from Table 1 with $\kappa = 0.2$ and $\gamma = 3$ (left panel); $\mathcal{N}(\hat{\mu}_r, \hat{\Sigma}_r)$ as above vs. $g_r^{KT}(\mu_r, \Sigma_r)$ with $\mu_r := 0.0115$ and Σ_r determined by CPT-style adjustment, corresponding to “Model, g^{KT} ” from Table 3 with any value of γ (right panel)).

that transforms the sample, or the observed probability density function $h_r(r)$ into what we argue is an estimate of a population, or true probability density $g_r(r)$ is available on the left panel of Figure 1.

Additional indicators: Table 2 provides some additional macroeconomic and financial data on the United States after World War II, together with the corresponding characteristics for the model specifications emphasized in the previous paragraph. One additional indicator is the correlation between the returns and consumption growth, which, being notably below 1, compares adequately to its empirical counterpart (in view of the fact that we are considering just a single risky asset).

The two other added indicators are the Sharpe ratio and the Hansen-Jagannathan bound. The first is defined as

$$\text{SR} := \frac{\text{E}[R] - R_0}{\sqrt{\text{V}[R]}} \leq \frac{\sqrt{\text{V}[M]}}{\text{E}[M]} =: \text{HJ}, \quad (4)$$

where the bound on the right-hand side defines the second one. The former quantity reflects the expected risk-return trade-off. The latter was formulated by Hansen and Jagannathan (1991) basing on the statistical properties of a valid SDF; it serves an im-

deviations of, respectively, consumption and dividend growth over 600,000 years of simulated data featuring rare booms and disasters in Tsai and Wachter (2016); it is also consistent with the finding in Malloy et al. (2009) that consumption growth of stockholders exhibits a sensitivity of factor about 3 to 4 relatively to aggregate consumption growth.

Table 2: Calibration Results, with Additional Details

	κ	γ	$E[R_0]$	$\sqrt{V[R_0]}$	$E[R - R_0]$	$\sqrt{V[R]}$	$E[\Delta c]$	$\sqrt{V[\Delta c]}$	$E[\Delta d]$	$\sqrt{V[\Delta d]}$	$\rho(r, \Delta c)$	SR	HJ
Empirical, h	n/a	n/a	0.88	1.26	8.23	17.05	1.94	0.99	2.43	4.23	0.18	0.23	n/a
Model, g	0.2	1	0.36	n/a	8.78	29.80	8.07	29.73	8.07	29.73	1.00	0.14	0.15
Model, g	0.2	2	0.34	n/a	8.79	29.80	2.87	14.64	2.87	14.64	0.99	0.14	0.15
Model, g	0.2	3	0.33	n/a	8.80	29.80	1.66	9.73	1.66	9.73	0.92	0.14	0.15
Model, g	0.2	4	0.33	n/a	8.80	29.80	1.16	7.29	1.16	7.29	0.85	0.14	0.15
Model, g	0.2	5	0.35	n/a	8.80	29.80	0.89	5.83	0.89	5.83	0.78	0.14	0.15
Model, g	0.1	3	-6.36	n/a	16.28	40.36	0.89	13.10	0.89	13.10	0.92	0.19	0.20
Model, g	0.2	3	0.33	n/a	8.80	29.80	1.66	9.73	1.66	9.73	0.92	0.14	0.15
Model, g	0.3	3	2.64	n/a	6.39	25.44	1.92	8.32	1.92	8.32	0.92	0.12	0.12
Model, g	0.4	3	3.79	n/a	5.22	23.01	2.05	7.53	2.05	7.53	0.92	0.11	0.11
Model, g	0.5	3	4.46	n/a	4.54	21.47	2.12	7.03	2.12	7.03	0.92	0.10	0.11

Notes: The columns present the mutual information parameter, the coefficient of relative risk-aversion; the real risk-free and market (excess) returns, the real per capita consumption growth as well as the real dividend growth with their statistical moments; the Sharpe Ratio denoted as “SR”, and the Hansen-Jagannathan bound denoted as “HJ”. The “Empirical” row: the data sample is U.S. 1948:Q1–2014:Q4, at quarterly frequency (see Appendix §C for a description of data sources). The “Model” rows: the CRRA utility, $\beta = 0.99$, and the probability density of returns (h or g) as the models’ only inputs (see text for a detailed description). The measurement units are nats for the mutual information distance, and percentage points converted into annualized terms for the economic variables.

portant role in the diagnostics of the asset-pricing models, often posing an impenetrable barrier for a theoretical Sharpe ratio to catch up with its empirical measurements.

A Sharpe ratio around 0.14 is below the empirically observed value of 0.23, but is closer than what is the case for the baseline “Model, h ” of Table 1, where the Sharpe ratio and the Hansen-Jagannathan bound are both equal to 0.08 (not shown in the tables). Of course, the improved performance of the latest model is natural, looking at the adjustment through the lens of Hansen-Jagannathan bound. Raising the variance of returns, by model’s general equilibrium restrictions, simultaneously increases the other relevant variances, including that of the SDF, which directly leads to an increase in the numerator of the HJ bound.⁷

Interpretation: Our preferred interpretation of these calibration results is as follows. The model we use is deliberately primitive, even crude perhaps, although truly respecting the general equilibrium discipline. It takes as inputs only the probability distribution of the risky asset’s returns, a commonly accepted level of the subjective discount factor β , and a RRA coefficient γ that is allowed to take one of the plausible values. The probability distributions we consider are the empirically observed distribution of returns (which we assume to take a parametric log-Normal form with parameters estimated from the available sample data) and the inferred “true” (population) probability distribution of returns. Anyone who takes the (joint) empirical probability distribution at face value (as “naïvely” does an econometrician, who uses the sample variance $\hat{\Sigma}_r$ and the sample expected return $E^h[R_{t+1}]$) finds it inconsistent with the simple model outlined above, in particular he ends up being puzzled by the “high” premium on risky assets and the “low” return on riskless assets.

On a closer look, however, this is just an illusion. More sophisticated agents (“professional” investors) recognize the fact that the stochastic environment represented by the population probability distribution is more complex, more risky than what an empirical probability distribution makes of it. Inferring their beliefs about the unobserved complex “true” distribution indirectly allows to rationalize the actual consumption and investment choices that lead to the observed (and erroneously perceived as “high”) premium on the risky assets and the (correspondingly, “low”) return on the risk-free assets. Even though investors may disagree with the simplistic empirical probability distributions and consider them inappropriate for investment decisions, the observed market levels of risk-free and risky returns do not look off from their perspective.

⁷Strictly speaking, matching Sharpe ratios is disputable as a relevant yardstick, because in contrast to traded assets, “disagreeing” with a Sharpe ratio does not provide investors with incentives to initiate the opposite investment/trade position.

2.2 Experiments

Taking a different tack, we are going to use the stylized facts gathered in numerous laboratory experiments, though not directly, but through the lens of the prospect theory of Kahneman and Tversky that summarizes them very accurately and compactly.

The prospect theory (Kahneman and Tversky, 1979) together with its refinement, the cumulative prospect theory (Tversky and Kahneman, 1992), applies to choice under risk and uncertainty and captures vast experimental evidence collected on such choices, offering robust empirical predictions. In particular, it is consistent with Allais' (1953) experimental results that challenge the conventional expected utility theory of choice. However, as the authors themselves emphasize, the (cumulative) prospect theory is not normative/prescriptive but rather is positive/descriptive, purely phenomenological. Its modern interpretation is that this is not a preference theory but a theory of "default actions" that are "appropriate to maximize experienced utility only on average" (Bossaerts et al., 2008), and which are actively used by less experienced agents (see Fox et al., 1996; List, 2004; also see List and Haigh, 2005).

The (cumulative) prospect theory stipulates that (after the initial framing stage) the overall value (or utility) of an uncertain "prospect" is a sum of values of the outcomes each multiplied by the decision weight. The outcome is defined with respect to a reference point, the value function potentially differs for positive ("concave for gains") and negative ("convex and steeper for losses") deviations from the reference point, and the decision weights possibly do not coincide with the presented probabilities (they are not even necessarily being additive to 1). It is the decision weight aspect that we are chiefly interested in. In the leading formulation of the theory, the decision weight function is defined as a non-linear transformation of the cumulative probability distribution function of the outcome, so that it overweights the small probabilities and underweights the moderate and large ones; hence the term "distortions" is also used commonly (see Figure 2 for a primer).

The value function-related prescriptions of the prospect theory can be reconciled with the felicity measurement side of our workhorse model by designating as a reference point wealth $W_t = 0$ in the value function and equivalently consumption $C_t = 0$ in the utility function (closing down the issue of asymmetric treatment of gains and losses due to the non-negativity of W_t and C_t).

Next, assuming that the degree of adjustment measured during the experiments of Kahneman and Tversky indeed reveals the optimal (in constrained, second-best sense)⁸ behavior, the effect of probability weighting allows us to estimate Σ_r and, in turn, indirectly calibrate κ .

We take the parametric measurements of the (cumulative) prospect theory's probability weighting function from Camerer and Ho (1994), who use data from nine different experimental studies to estimate the parameters we are interested in. The functional form

⁸Cf. Herbert Simon's "ecological rationality" and Gerd Gigerenzer's "shortcut heuristics".

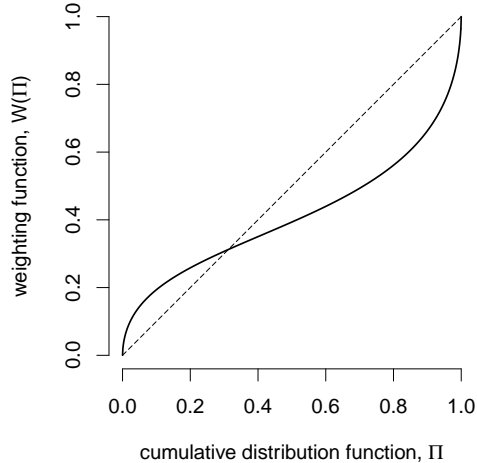


Figure 2: Cumulative prospect theory’s probability weighting function primer.

they use is

$$W(\Pi) := \frac{\Pi^\eta}{(\Pi^\eta + (1 - \Pi)^\eta)^{1/\eta}},$$

where $\Pi(\cdot)$ is some cumulative distribution function, $W(\cdot)$ is the cumulative probability weight for a given partition of the probability space, and $\eta \leq 1$ is a parameter, with $\eta = 1$ corresponding to the standard linear probability weighting. Camerer and Ho’s estimate of η is 0.56.

Application: Now, we apply the prospect theory-style probability weighting adjustment using the functional form and parameter estimate of Camerer and Ho (1994). As a reference value, with respect to which the adjustment procedure determines what probability quantile a particular realization belongs to, we use μ_r , the mean of the distribution we end up with.⁹ The right panel of Figure 1 shows the end product of the probability weighting adjustments implied by the prospect theory of Kahneman and Tversky, depicted as $g_r^{KT}(r)$, along with the sample distribution, $h_r(r)$.

Visually, the similarity between the two putative true distributions $g_r(\cdot)$ and $g_r^{KT}(\cdot)$ is striking, their location and scale parameters seem well attuned. A quantitative comparison is provided by Table 3, which presents the results of an exercise analogous to the one conducted for “Model, g ” with $\kappa = 0.2$ and $\gamma = 3$ in Tables 1 and 2, but this time using the distribution $g_r^{KT}(\cdot)$ that is produced following the stipulations of the prospect theory. For $\gamma = 3$, the procedure of Kahneman and Tversky turns out the results that roughly match those we got in the calibrations of §2.1.¹⁰

⁹Which is obtained by finding such a value for μ_r that would ensure the expected value of the return is preserved, $E^g[R_{t+1}] = E^h[R_{t+1}]$.

¹⁰Historically, this is not the first paper that refers to the prospect theory-style probability weighting in rationalizing the observed equity premium and related empirical phenomena. The others include an early attempt by Epstein and Zin (1990), as well as Routledge and Zin (2010), De Giorgi and Legg (2012), also notwithstanding the related contribution of Barberis and Huang (2008).

Table 3: Calibration Results for Kahneman-Tversky Approach, with Additional Details

	κ	γ	$E[R_0]$	$\sqrt{V[R_0]}$	$E[R - R_0]$	$\sqrt{V[R]}$	$E[\Delta c]$	$\sqrt{V[\Delta c]}$	$E[\Delta d]$	$\sqrt{V[\Delta d]}$	$\rho(r, \Delta c)$	SR	HJ
Empirical, h	n/a	n/a	0.88	1.26	8.23	17.05	1.94	0.99	2.43	4.23	0.18	0.23	n/a
Model, g	0.2	3	0.33	n/a	8.80	29.80	1.66	9.73	1.66	9.73	0.92	0.14	0.15
Model, g^{KT}	n/a	1	0.63	n/a	8.47	29.19	8.05	29.11	8.05	29.11	1.00	0.14	0.14
Model, g^{KT}	n/a	2	0.70	n/a	8.43	29.10	2.92	14.34	2.92	14.34	1.00	0.14	0.14
Model, g^{KT}	n/a	3	0.71	n/a	8.42	29.09	1.71	9.53	1.71	9.53	0.90	0.14	0.14
Model, g^{KT}	n/a	4	0.70	n/a	8.41	29.06	1.19	7.13	1.19	7.13	0.77	0.14	4.31
Model, g^{KT}	n/a	5	0.75	n/a	8.39	29.03	0.92	5.69	0.92	5.69	0.63	0.14	3193.06

Notes: The columns present the mutual information parameter, the coefficient of relative risk-aversion; the real risk-free and market (excess) returns, the real per capita consumption growth as well as the real dividend growth with their statistical moments; the Sharpe Ratio denoted as “SR”, and the Hansen-Jagannathan bound denoted as “HJ”. The “Empirical” row: the data sample is U.S. 1948:Q1–2014:Q4, at quarterly frequency (see Appendix §C for a description of data sources). The “Model” rows: the CRRA utility, $\beta = 0.99$, and the probability density of returns (h , g or g^{KT}) as the models’ only inputs (see text for a detailed description). The measurement units are nats for the mutual information distance, and percentage points converted into annualized terms for the economic variables.

More interpretation: The consistency between the results in Tables 1–2 and Table 3 suggests the same putative population distribution at the mutual information distance of $\kappa = 0.2$ nats from the sample distribution. Moreover, later we will argue that these two sets of empirical results are not just related, but represent different sides of the same coin.

2.3 Additional empirical results

The existing literature has accumulated some further empirical results associated with probability weighting transformations as in the (cumulative) prospect theory of Kahneman and Tversky.

Take the risk-aversion/IMPLIED volatility “smile” pattern that is observed empirically in the options segment of financial markets. Fundamentally, this “non-monotone pricing kernel puzzle” is a phenomenon that proves to be challenging for conventional theoretical models to convincingly address (e.g., an overview can be found in Ziegler, 2007). However, the fact that the sample variance-covariance matrix of returns is smaller than its population counterpart, and one way to counteract this is to adjust it upwards—for instance, using probability weighting as stipulated by the prospect theory of Kahneman and Tversky—has direct relevance to the above phenomenon in finance.

Specifically, Kliger and Levy (2009) show that prospect theory-style probability weighting helps explain observed options prices. In contrast to the parametric approach taken in the previous paper, Polkovnichenko and Zhao (2013) use an agnostic non-parametric procedure to estimate an empirical SDF as well as a probability density function from options data, and then (conditionally on the utility function specification assumptions) extract empirical probability weighting functions, which turn out to deviate from linearity and to be in line with the probability adjustments presented above.

Another such regularity is the “portfolio concentration/underdiversification puzzle” (see Blume and Friend (1975), Statman (1987), Kelly (1995) for individual investors’ domestic portfolios; French and Poterba (1991) for country-level international portfolios). Again, apparent “underdiversification” may be an artifact of using a smaller, sample variance in place of a larger, population one: with covariances held fixed, employing the former mechanically amplifies the correlation coefficients. As a result, investment portfolios would look less diversified than they actually are.

Quantitatively speaking, the work of Polkovnichenko (2005) shows that these empirical observations can also be explained by the probability weighting of the type stipulated by Kahneman and Tversky’s theory.

3 Theory

3.1 Conceptual approach

Generally speaking, the problem we are dealing with is about learning, inference and decision in a situation of scarce data. For instance, small samples, rare events, rich stochastic dynamics all present a challenge of this kind.

A necessary step in resolving this problem is to specify what we know and what we do not know, and then to express our ignorance in a rigorous way. The “principle of indifference” offers a disciplined approach to tackling such ignorance, stipulating that equivalent states of knowledge should be assigned equivalent probabilities. In practice this principle is usually implemented through the method of invariance or transformation groups (Jeffreys, 1939; 1946) and the method of maximum entropy (Jaynes, 1957a; 1957b). (Rissanen, 1983, proposed a unifying understanding of these two methods based on his “minimum description length principle”.)

More concretely, taking a Bayesian perspective on probability theory, the above approach boils down to specifying very carefully a non-informative prior probability distribution that possesses the desired invariance properties. Ultimately, combining the information contained in the prior distribution with the information provided by sample data (if any) defines a stochastic environment that is relevant for a particular decision choice.

A few points deserve elaboration. First, the prior distributions we use are characterized by invariance to transformations. The method of transformation groups ensures that the specification of a given problem (where specification also includes any prior probabilistic information we might have) is unaffected by equivalent transformations applied to it. In a one-dimensional case, as here, this approach coincides with the Jeffreys’ method of assigning prior distributions that would be invariant to reparameterization. For instance, it delivers translation invariance for a location parameter (mean in the case of Gaussian distribution) and measurement-units invariance for a scale parameter (variance, respectively).

Second, the above prior distributions are characterized by maximal entropy within their corresponding (in some sense) classes of probability distributions. The method of maximum entropy ensures that the prior probability distributions considered reflect no more and no less than the information at hand. This is achieved by maximizing a certain measure of uncertainty (i.e., our lack of knowledge) taking as given the available information, explicitly specified; or, equivalently, by minimizing the amount of “residual” information.

Lastly, the ignorance prior distributions we thus obtain turn out to be conjugate priors, which means that their posterior updates are preserving the distributional form.¹¹

¹¹Roughly speaking, this is due to the fact that we are using likelihoods from the exponential family of probability distributions. Those admit conjugate prior distributions that also belong to the exponential

3.2 Lotteries in laboratory experiments

Consider a simple lottery such that its outcome is represented by a binary random variable $x \in \{0, 1\}$ with a probability of successful outcome 1 being equal to $p \in [0, 1]$. The corresponding Bernoulli probability mass function defines the likelihood of possible outcomes:

$$L(x|p) := x^p(1-x)^{1-p}. \quad (5)$$

Before facing the lottery mechanism (say, before entering the experimental laboratory) and learning anything about the characteristics of the lottery, which in this primitive example is captured exhaustively by the value of parameter p , the agent formulates a prior distribution that accurately reflects his (lack of) knowledge about the above parameter. The Jeffreys' method prescribes a prior proportional to the square root of the determinant of the Fisher information matrix¹², which in our case yields:

$$\pi(p) \propto \sqrt{E^L \left[\left(\frac{\partial}{\partial p} \ln L(x|p) \right)^2 \right]} = \frac{1}{p^{1/2}(1-p)^{1/2}}. \quad (6)$$

This constitutes the kernel of a Beta probability distribution $\mathcal{B}(\alpha, \beta)$ with parameters $\alpha = 1/2$ and $\beta = 1/2$. The corresponding symmetric U-shaped distribution (as plotted in Figure 3) has most of its probability mass concentrated in the poles on each end of the $[0, 1]$ interval. This is different from the naïve non-informative prior that implies p uniformly distributed on the real $[0, 1]$ interval (and which is obtained from the so-called Bayes–Laplace prior $\mathcal{B}(1, 1)$). Instead, parameterization invariance implies that a transformed parameter $\arcsin(p^{1/2})$ has a uniform distribution on the $[0, \pi/2]$ interval.

Thus, we obtained the Beta–Bernoulli conjugate system. Intuitively, using the trick that any conjugate prior can be viewed as the posterior for a pseudo data set, as well as a standard interpretation of the Beta distribution's parameters α and β as pseudo-observations or pseudocounts of Bernoulli trials, values $\alpha = 1/2$ and $\beta = 1/2$ represent 1 pseudo-observation equally split between a success and a failure (while the naïve Beta–Laplace prior represents 1 successful and 1 unsuccessful pseudocounts).

Turning to the maximum entropy approach, the Beta distribution maximizes entropy on the $[0, 1]$ interval subject to two symmetric constraints on the logarithm of the geometric mean of the random variable p : $E^\pi[\ln(p)] =: \psi(\alpha) - \psi(\alpha + \beta)$ and $E^\pi[\ln(1-p)] =: \psi(\beta) - \psi(\alpha + \beta)$, where $\psi(\cdot)$ is the digamma function. Thus, the values of α and β above may be alternatively interpreted as a specific restriction on the geometric mean value in entropy maximization.¹³

family. In turn, the exponential family distributions arise naturally as a result of entropy maximization. Moreover, for the distributions in the exponential family, the Jeffreys priors are optimal in the “minimum description length” sense.

¹²This ensures the parameterization invariance rather “mechanically” by preempting the appearance of terms that arise from reparameterization using the standard change of variables procedure.

¹³For the sake of completeness, it is $1/4$ for the probability of success, and $1/4$ for the probability of

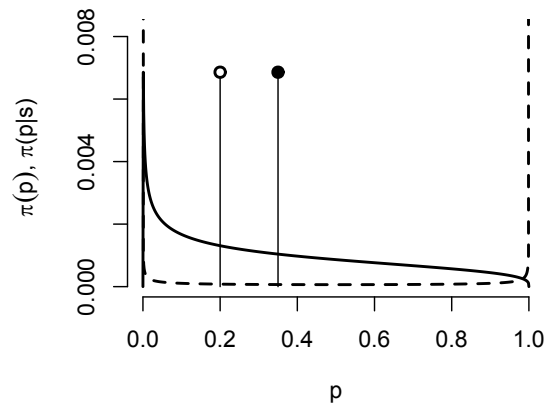


Figure 3: Posterior (solid) vs. prior (dashed) probability densities (parameterizations used are as follows: prior $\pi(p)$ is $\mathcal{B}(\alpha, \beta)$ with $\alpha := 1/2$ and $\beta := 1/2$ (re-normalized to fit the plot), posterior $\pi(p|s)$ is $\mathcal{B}(\alpha, \beta)$ with $\alpha := 1/2 + 1 \times 0.20$ and $\beta := 1/2 + 1 \times (1 - 0.2)$, announced value $s := 0.20$ (spike with empty bullet on top), posterior mean $E[p|s] := 0.35$ (spike with filled bullet on top)).

Now (having entered the lab), the agent learns something about the lottery: he is told that the probability of success is s .

Note that there is no sample data yet as the agent still has not experienced any lottery trials; he only has some non-sample information (possibly vague or non-credible) that can be used to modify the prior distribution: e.g., to bias it one way or another, to make it sharper or more diffuse, etc. This is in the spirit of “no-data decision problem” due to Chernoff and Moses (1959). In such cases, the prior to some extent determines the conclusion (inference or decision).

Then the agent updates the initial prior with the sufficient statistic that corresponds to $p = s$ using the standard mechanism of Bayesian updating, that is through adding new (pseudo-) observations. For the Bernoulli distribution, the sufficient statistic is the share of successes, so for the announced probability s the value is going to be s itself, multiplied by the number of pseudo-observations $n \in \mathbb{N}$.

Consequently, the posterior distribution (the updated prior, really) is

$$\pi(p|s) = \frac{p^{\alpha+ns-1}(1-p)^{\beta+n(1-s)-1}}{B(\alpha+ns, \beta+n(1-s))}, \quad (7)$$

which, due to conjugacy, is also of the Beta form (where $B(\cdot, \cdot)$ denotes the Beta function).

If the agent has full confidence in the lottery mechanism (i.e., in the lottery organizers, or in the device itself), then the number of pseudo-observations is infinite, the prior gets washed-out, and the posterior distribution becomes concentrated at the announced value s . This is the case that is conventionally considered in the literature. Thus, the sharpness of the updated prior depends on n , which effectively reflects the relative degree of confidence in the lottery mechanism.

However, there is no a priori rationale to rule out the less-than-full confidence case of $n < \infty$. If there is no reason to weight one type of prior information higher than the other, by the principle of indifference between sources of pseudo-observations, $n := 1$ is the value that puts two sources on equal footing and treats them equivalently.¹⁴ Intuitively, such an approach can be interpreted as extending the existing 1 pseudo-observation from the initial prior with 1 additional pseudo-observation from the external source (which in turn is itself composed of s pseudocounts of successes and $(1-s)$ pseudocounts of failures).

The resulting posterior distribution of p , defined by equation (7) with $n = 1$, is not centered at the announced value s but instead is biased towards the extremities of the $[0, 1]$ interval. In particular, the mean of the posterior distribution is a ratio¹⁵

$$E[p|s] := E^{\pi(p|s)}[p] = \frac{\alpha + ns}{\alpha + ns + \beta + n(1-s)} = \frac{\alpha + s}{\alpha + \beta + 1}. \quad (8)$$

failure — they do not have to add up to 1 (in the Bayes–Laplace case, the corresponding values are $1/e$ and $1/e$).

¹⁴The argument here is heuristic, but it can be stated more formally, e.g. probabilistically by the method of maximum entropy that stipulates a uniform distribution of discrete probability weights, or by the invariance utilizing the symmetry of information sources.

¹⁵Cf. the usual number-of-observations–weighted average in canonical Bayesian updating.

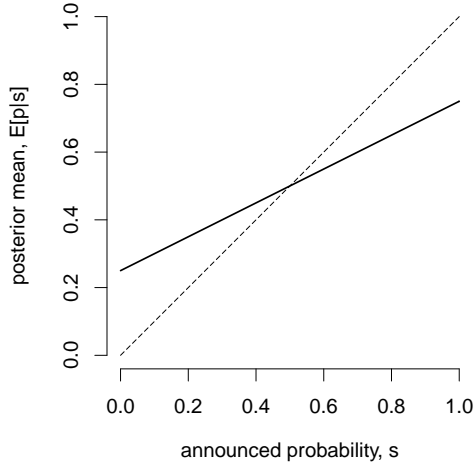


Figure 4: Posterior-mean (vertical axis) vs. announced (horizontal axis) probabilities (parameterization used is $E[p|s]$ from equation (8) with $\alpha := 1/2$, $\beta := 1/2$ as well as $n := 1$).

For example, an announced value of s that is equal to 0.20 results in a posterior mean of p being 0.35 (see Figure 3 for the posterior shape and mean). This is reminiscent of the probability weighting discussed in part §2.

Indeed, juxtaposing the full range of announced probabilities s against the corresponding posterior means of p (as presented in Figure 4) is broadly consistent with the pattern of probability weights we have seen in part §2 (i.e., in Figure 2). Specifically, low probabilities are increased, high probabilities are decreased, and the emerging distortions are larger for more extreme probabilities.

In the presented derivations and the resulting plot, the probability distortions exhibit a linear pattern; they have a fixed, or inflection, point at $p = 1/2$; and distorted probabilities in the corners of $p = 0$ and $p = 1$ do not converge continuously, or “paste smoothly”, with the original probabilities (so a degenerate, non-stochastic lottery becomes stochastic). These features contradict the common view on the relevant stylized facts (as given in, say, Kahneman and Tversky, 1979, 1992; or in Camerer and Ho, 1994): an inverse-S-shaped pattern; a fixed point around $p = 1/3$; the “smooth pasting” at 0 and 1 (e.g., see Figure 2). However, more recent estimates that use less restrictive non-parametric methods and focus on representative samples of participants are in fact more supportive of the linear pattern of distortions as well as the inflection at the midpoint (see Fehr-Duda and Epper, 2012, for some evidence). Moreover, degenerate lotteries are rarely tested in the experiments, hence the “smooth pasting” result may be an artifact of the parametric estimation and fitting procedure, and, strictly speaking, warrants a confirmation in a separate, narrowly-focused study.

Finally, with the updated prior as above, the agent is optimally equipped for making the decision regarding the proposed lottery (to choose the lottery by, say, pressing a key in the lab). Moreover, the agent is also equipped for learning from his decisions, if more

than one trial of the same lottery is run; indeed, given such an ergodic and stationary environment, the posterior distribution will converge to a Dirac delta function at p and eventually he will learn the true lottery parameter perfectly.

The benefits of the presented indifference-based approach to tackling ignorance about important economic parameters are as follows. Of course, the main motivation is not our toy example of lotteries in laboratory experiments, but rather learning some probabilities of interest empirically, in real time/online. A good practical example would be making decisions in stochastic environments on a regular basis, frequently encountering fresh lotteries with unknown parameters in the process (to buy some durable consumption item or not? to work on some task for another hour or to rest instead? to invest some unexpected income or to give to charity?). In particular, the updated prior above is beneficial when learning in an undersampled setting: in environments characterized by “rare events” and under conditions of “small samples”. This is achieved by reducing the reliance on observed data and allowing, in an optimal way, for unseen contingencies.¹⁶

Refinements: The above derivations and results are our preferred way of thinking about this problem. Nevertheless, some further refinements of the general approach that would rationalize certain potentially desirable stylized empirical facts are presented below.

First, the matching of announced and posterior (i.e., of original and distorted) probabilities at the corners $p = 0$ and $p = 1$ would be ensured for an agent who uses, instead of the Jeffreys prior $\mathcal{B}(1/2, 1/2)$, a mixture prior that combines the Jeffreys one with the so-called Haldane prior $\mathcal{B}(0, 0)$. The “distribution” that corresponds to an improper Haldane prior is essentially a convolution of two Dirac delta functions on each end of the $[0, 1]$ interval. Instead of a parameterization invariance, Haldane’s approach is motivated by the fact that it effectively imposes the posterior mean of the distribution of the parameter in question to coincide with its maximum-likelihood point estimate. Still, it does imply a uniform distribution for a particular parameterization: in terms of the logarithm of the odds ratio, $\ln(p/(1-p))$. Intuitively, it can be interpreted as representing 0 initial pseudo-observations.¹⁷ Standard Bayesian updating in this case produces the posterior

$$\pi^M(p|s) = \sum_{m=1}^M \pi(m|s) \times \frac{p^{\alpha_m + ns - 1} (1-p)^{\beta_m + n(1-s) - 1}}{\mathcal{B}(\alpha_m + ns, \beta_m + n(1-s))}, \quad (9)$$

with

$$\pi(m|s) := \frac{\pi(m) \mathcal{B}(\alpha_m + ns, \beta_m + n(1-s)) / \mathcal{B}(\alpha_m, \beta_m)}{\sum_k \pi(k) \mathcal{B}(\alpha_k + ns, \beta_k + n(1-s)) / \mathcal{B}(\alpha_k, \beta_k)},$$

¹⁶Lastly, it may sound implausible that ordinary people participating in laboratory experiments perform the above Bayesian updating-related computations. Indeed, our understanding is that they are not consciously doing these calculations, but rather acquire evolutionarily and/or learn experientially the optimal behavior/policy rules by the usual method of trial and error.

¹⁷For completeness, the maximum entropy interpretation here implies a geometric mean value restriction of 0 both for the probability of success as well as for the probability of failure.

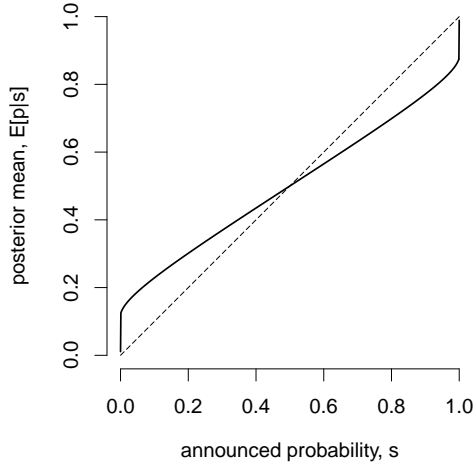


Figure 5: Posterior-mean (vertical) vs. announced (horizontal) probabilities, refinement (parameterization used is $E[p|s]$ utilizing equation (9) with $\alpha_m, \beta_m \in \{0.50, 0.40, 0.31, 0.23, 0.16, 0.10, 0.06, 0.03, 0.01, 0.00\}$ and $\pi_m \in \{0.01, 0.01, 0.01, 0.01, 0.01, 0.01, 0.01, 0.01, 0.01, 0.01, 0.9\}$ as well as $n := 1$).

where $\pi(m)$ is the initial prior weight of the mixture component m . For a high enough $\pi(m)$ corresponding to the Haldane prior, and with $M := 2$, we get the corner probabilities being matched, while the interior ones are unaffected because a dogmatic Haldane prior has zero posterior weight in those regions (hence, Figure 4 is unchanged, bar its two end points). However, those corner probabilities do not “paste smoothly”, there are discontinuous jumps on each end.

Second, the “smooth pasting” of the announced and posterior probabilities would be achieved by expanding the set of the mixture prior’s components with the prior distributions located (in the hyperparameter space) between the Jeffreys’ and the Haldane’s one. Using the same Bayesian updating formula (9), with $\{\alpha_m, \beta_m\}$ spread between $\{0, 0\}$ and $\{1/2, 1/2\}$, as well as with initial prior weights $\pi(m)$ equally allocated between mixture components (except for the disproportionately high weight on the Haldane’s component in the corner), we get the corner probabilities “smoothly pasted” together with an inverse-S-shaped pattern of probability distortions (e.g., see Figure 5).

3.3 Assets in financial markets

As to the more practically interesting case with continuous distributions, we can get a sense of what would happen by taking the part’s §3.2 Bernoulli distribution to the limit. According to the classical de Moivre–Laplace theorem, as n grows large, the distribution of the number of successes in the Bernoulli trials nx will converge to the Normal distribution with mean np and variance $np(1 - p)$. Taking, without loss of generality, $p = s > 1/2$, by the derivations above the posterior mean of p would be biased towards $1/2$, which unambiguously increases the corresponding Normal distribution’s variance due to the

latter's concavity in p .¹⁸

More realistically, we now focus on a non-ergodic environment with continuous, unbounded distributions. Thus, even a large number of random variables' realizations would never allow perfect learning of the (ever-evolving) true parameters, and the prior could not possibly become dominated by the sample observations. (In terms of lotteries from part §3.2, the non-ergodicity would correspond to running non-repeated lottery trials).

Consider stock price returns whose logarithms r_t are Normally distributed with mean μ_t and variance σ_t^2 . The corresponding Gaussian probability density function defines the likelihood of possible returns:

$$L(r_t|\mu_t, \sigma_t) := \frac{1}{\sqrt{(2\pi\sigma_t^2)}} \exp\left(-\frac{1}{2\sigma_t^2}(r_t - \mu_t)^2\right). \quad (10)$$

Here, the parameters themselves are (unobserved) random variables. To obtain the corresponding ignorance prior invariant under reparameterization, the Jeffreys' method yields (treating mean and variance parameters separably, rather than as a single vector¹⁹):

$$\pi(\sigma_t^2|\cdot) \propto \sqrt{\mathbb{E}^L \left[\left(\frac{\partial}{\partial \sigma_t^2} \ln L(r_t|\mu_t, \sigma_t^2) \right)^2 \right]} = \frac{1}{\sigma_t^2}, \quad (11)$$

$$\pi(\mu_t|\cdot) \propto \sqrt{\mathbb{E}^L \left[\left(\frac{\partial}{\partial \mu_t} \ln L(r_t|\mu_t, \sigma_t^2) \right)^2 \right]} = \frac{1}{\sigma_t}. \quad (12)$$

These are improper priors, with the one corresponding to the mean found to be conditional on the variance. Indeed, we can always write $\pi(\mu_t, \sigma_t^2) = \pi(\mu_t|\sigma_t^2)\pi(\sigma_t^2)$.

It will prove convenient to assume that the variance σ_t^2 follows an Inverse-Gamma distribution²⁰ and the mean μ_t is Normally distributed (conditionally on the variance). In line with this, the above equations (11) and (12) can be viewed as limiting cases of an Inverse-Gamma probability distribution $\mathcal{IG}(\alpha_0, \beta_0)$ with parameters $\alpha_0 = 0$ and $\beta_0 = 0$, as well as a (conditional) Gaussian distribution $\mathcal{N}(\ell_0, \sigma_t^2/\lambda_0)$ with parameters $\ell_0 \in \mathbb{R}$ (say, $\ell_0 = 0$) and $\lambda_0 \rightarrow 0$. Thus, parameterization invariance implies that a transformed scale parameter $\ln \sigma_t^2$ has a uniform distribution on the real line, leaving the scale to be multiplication-invariant; while a location parameter μ_t has a (conditionally) uniform distribution on the real line, being addition-invariant.

Thus we obtained the Normal-Inverse-Gamma-Normal conjugate system. Intuitively, the Normal prior distribution's parameter λ_0 and Inverse-Gamma's parameter α_0 can be interpreted as pseudo-observations of the Normal realizations of r_t — none in our case.

¹⁸Provided we ignore the effect of learning from different trials (say, by treating each partition of the probability space/every cell of the domain as a genuinely new lottery).

¹⁹There is still a lot of disagreement in the existing Bayesian literature on how to deal with invariance for location-scale parameters in the exponential-family distributions, so some personal judgement is necessary here. Relevant sources include DeGroot (1970), Kass and Wassermann (1996), Minka (2001), Gelman et al. (2004), Murphy (2012).

²⁰Here, $\mathcal{IG}(h, s)$ with shape h and inverse scale s is parameterized as $f(x|h, s) = \frac{s^h}{\Gamma(h)} x^{-h-1} \exp(-\frac{s}{x})$.

By the maximum entropy approach, the Inverse-Gamma distribution maximizes entropy on the $(0, \infty)$ interval subject to constraints on the arithmetic mean and on the logarithm of the geometric mean of the random variable $1/\sigma_t^2$: $E^\pi[1/\sigma_t^2] =: \alpha_0\beta_0$ and $E^\pi[\ln(1/\sigma_t^2)] =: \psi(\alpha_0) + \ln \beta_0$; while the Normal distribution maximizes entropy on the $(-\infty, \infty)$ interval subject to constraints on the (arithmetic) mean and on the variance of the random variable μ_t : $E^\pi[\mu_t] =: \ell_0$ and $E^\pi[(\mu_t - \ell_0)^2] =: \sigma_t^2/\lambda_0$. Thus, in general the values of α_0 and β_0 above may also be interpreted as specific restrictions on the mean (but not point) values of σ_t^2 ; while the values of ℓ_0 and λ_0 on the first moment and the second (central) moment of μ_t in entropy maximizations.²¹

To start with, consider an ergodic environment where σ_t^2 and μ_t are realized at time $t = 0$ and remain constant thereafter. In such situation, whenever new observations on returns r_t are realized, the standard mechanism of Bayesian updating stipulates:

$$\alpha_t := \alpha_0 + \frac{1}{2}t \quad (13)$$

$$\beta_t := \beta_0 + \frac{1}{2} \left(\frac{\lambda_0 t}{\lambda_0 + t} (\bar{r}_t - \ell_0)^2 + \sum_{i=1}^t (r_i - \bar{r}_t)^2 \right) \quad (14)$$

$$\ell_t := \frac{1}{\lambda_0 + t} (\lambda_0 \ell_0 + t \bar{r}_t) \quad (15)$$

$$\lambda_t := \lambda_0 + t \quad (16)$$

Consequently, the (marginal) posterior distributions for the variance and mean are, respectively,

$$\pi(\sigma_t^2 | r_t) = \frac{\beta_t^{\alpha_t}}{\Gamma(\alpha_t)} (\sigma_t^2)^{-\alpha_t-1} \exp\left(-\frac{\beta_t}{\sigma_t^2}\right), \quad (17)$$

$$\pi(\mu_t | r_t) = \frac{\Gamma(\frac{2\alpha_t+1}{2})}{\Gamma(\frac{2\alpha_t}{2}) \sqrt{2\pi\beta_t/\lambda_t}} \left(1 + \frac{(\mu_t - \ell_t)^2}{2\beta_t/\lambda_t}\right)^{-\frac{2\alpha_t+1}{2}}, \quad (18)$$

being, due to conjugacy, of the Inverse Gamma and Student's t form (and where $\Gamma(\cdot)$ denotes the Gamma function). The latter distribution is born from the Normal that was conditional on variance, and is the non-standardized Student's t distribution²² with degrees-of-freedom $2\alpha_t$, location ℓ_t and scale $\beta_t/(\alpha_t\lambda_t)$; in short $\mathcal{T}_{2\alpha_t}(\ell_t, \beta_t/(\alpha_t\lambda_t))$.

In turn, the predictive posterior distribution for next period's returns is

$$\pi(r_{t+1} | \ell_t, \lambda_t, \alpha_t, \beta_t) = \frac{\Gamma(\frac{2\alpha_t+1}{2})}{\Gamma(\frac{2\alpha_t}{2}) \sqrt{2\pi\beta_t(1+1/\lambda_t)}} \left(1 + \frac{(r_t - \ell_t)^2}{2\beta_t(1+1/\lambda_t)}\right)^{-\frac{2\alpha_t+1}{2}}, \quad (19)$$

which is another Student's t , $\mathcal{T}_{2\alpha_t}(\ell_t, \beta_t(1+1/\lambda_t)/\alpha_t)$.²³

²¹That is, ∞ variance and a free/unrestricted mean, respectively.

²²Here, $\mathcal{T}_d(l, s)$ with degrees-of-freedom d , location l and scale s is parameterized as $f(x|d, l, s) = \frac{\Gamma(\frac{d+1}{2})}{\Gamma(\frac{d}{2}) \sqrt{d\pi s}} \left(1 + \frac{(x-l)^2}{ds}\right)^{-\frac{d+1}{2}}$.

²³The above derivations, which pertain to our ergodic case, are standard in Bayesian analysis, e.g. see Murphy (2012).

Now, consider a non-ergodic environment where σ_t^2 and μ_t may change (to anywhere within $\mathbb{R}_+ \times \mathbb{R}$) at any time t^* (where t^* is such that a binary random variable $\chi_t \in \{0, 1\}$ switches to 1 at $t = t^*$); then the parameters remain constant until the next change at some time t^{**} ; and so on. Assume that (i) a parameter change occurs with probability p . For an ease of exposition, also assume that (ii) probability p is non-stochastic and is common knowledge; and that (iii) the time of change is known ex post but not ex ante (the parameter values are not known though and must be learned). Such an evolution of the underlying parameter values can be understood as changes in (latent) macro-financial regimes.²⁴

While staying within the same regime (and foreseeing this), the learning dynamics follow the formulas (13)–(16) given above. So, by (17)–(18), the parameters would eventually be learned. By (19), the posterior Student's t distribution of returns would eventually converge to the true distribution (10). However, these eventual target distributions will never be fully achieved as long as probability of regime change p is above 0.

Because at the time of regime change t^* , the above learning recursion restarts. That is, formulas (13)–(19) are re-initialized with $t := 1$. Since initially nothing is known about the new regime, the prior parameter values are reset too. For the variance, to $\alpha_0 = \beta_0 = 0$, implying the logarithmic prior introduced earlier. For the mean, to $\lambda_0 \rightarrow 0$, which implies infinite variance and effectively no information about the location (in the same vein, ℓ_0 , the centrality parameter for the mean, remains unrestricted, but this is of no effect given $\lambda_0 \rightarrow 0$).

However, in the absence of empirical observations from the new regime, the above parameter values lead to improper prior distributions. Such priors can not be used for predictions and decisions that could account for the possibility of a future move to a new regime. To avoid this deficiency, we will bring in some additional information by updating these priors with the sufficient statistics that correspond to the regime immediately preceding the future one, i.e. with some statistics from the current regime.

Specifically, the agent updates the prior by adding the empirical mean $\bar{r}_t \approx \mu_t$ and (cumulative) variance $\sum_{i=1}^t (r_i - \bar{r}_t)^2 =: s_t^2 \approx t\sigma_t^2$ via, respectively, equations (15) and (14). The rationale would be the following: if there is no reason to favor μ_t and $\ln \sigma_t^2$ smaller or larger than what has been observed in the current regime, by the principle of indifference between the respective left and right orthants, if we must make a pick, then μ_t and $\ln \sigma_t^2$ are the values that treat two orphans equivalently.²⁵ Then, for proper prior distributions (in particular, for positive α and β), the minimally required value of t is 1.

²⁴Although we abstract away from it, the methods for inferring regimes and updating the probability of change are well-known, extensive reviews can be found in Ang and Timmermann (2012) or Hamilton (2016).

²⁵This argument is heuristic, so readers who are uncomfortable with our approach may prefer to view these updated prior distributions as reflecting the knowledge of the true parameter change process (conditionally on a change actually occurring, $\chi_t = 1$), which in this case is supposed to be a random walk of $\{\ln \sigma_t^2, \mu_t\}$ with increments uniformly distributed on \mathbb{R}^2 .

This value, similarly to setting $n := 1$ in part §3.2, puts two sources of prior information on equal footing (to the extent possible). Intuitively, such an approach can be interpreted as extending the existing 0 pseudo-observations from the initial, Jeffreys prior with 1 pseudo-observation from the external source, the previous regime.

Thus, conditionally on regime change, the predictive posterior distribution, defined by equation (19) with the described above α_1 , β_1 , λ_1 and ℓ_1 as well as $t := 1$, becomes a Student's t with 1 degree of freedom, $\mathcal{T}_1(\bar{r}_t, 2s_t^2)$, i.e. a Cauchy distribution.

Unconditionally, the resulting predictive posterior distribution for r_{t+1} becomes a mixture of two Student's t distributions:

$$\pi^M(r_{t+1}|\ell_t, \lambda_t, \alpha_t, \beta_t; \bar{r}_t, s_t^2) = (1-p) \times \pi(r_{t+1}|\ell_t, \lambda_t, \alpha_t, \beta_t) + p \times \pi(r_{t+1}|\bar{r}_t, 1, 1/2, s_t^2/2), \quad (20)$$

where $\pi(r_{t+1}|\ell_t, \lambda_t, \alpha_t, \beta_t)$ is $\mathcal{T}_{2\alpha_t}(\ell_t, \beta_t(1+1/\lambda_t)/\alpha_t)$ and $\pi(r_{t+1}|\bar{r}_t, 1, 1/2, s_t^2/2)$ is $\mathcal{T}_1(\bar{r}_t, 2s_t^2)$.

It is important to clarify some properties of the above mixture. The Student's t distribution is infinitely divisible but is not closed under convolution, so it is not stable (i.e., a sum of Student's t 's random variables is not a Student's t variable). As a result, characterizing the mixture distribution is not straightforward. The most challenging component is the Cauchy distribution, because all of its moments are undefined.²⁶ Nevertheless, in our case this does not pose any material difficulties, as stated in the following Proposition.

Proposition 1 (Mixture Moments). *The mixture distribution $\pi^M(\cdot)$ defined in equation (20) always possesses the first raw moment (mean); it also possesses the second central moment (variance) unless $\alpha_t < 1$.*

Proof. This is a result due to Theorem 3 in Berg and Vignat (2010), which shows that a sum of two independent Student's t random variables is an infinite convex combination of Student's t random variables, where each Student's t component has the number of degrees of freedom equal to the sum of degrees of freedom in the original two summands. After additions/subtractions, rescaling and rearrangement of the terms in (20), we can appeal to the above-mentioned Theorem; then, using the positivity and boundedness of its infinite series' weights, the existence result follows. □

The expected lifetime of a new regime (number of periods till the next parameter change, $E_t^{\pi(x)}[t^* - t]$) is $1/p$, which determines the corresponding $\alpha_{t^*} = 1/(2p)$, and thus the expected degrees-of-freedom in the non-Cauchy component $\mathcal{T}_{2\alpha_{t^*}}(\ell_t, \beta_t(1+1/\lambda_t)/\alpha_{t^*})$ of the Student's t mixture given by equation (20).

As an illustration, taking $p := 1/5$, which at the quarterly frequency adopted in

²⁶Though, to be fair, its cumulative distribution function is well defined, its quantile function is analytical (let alone trimming/winsorizing is always an option), enabling certain empirical usage of the Cauchy.

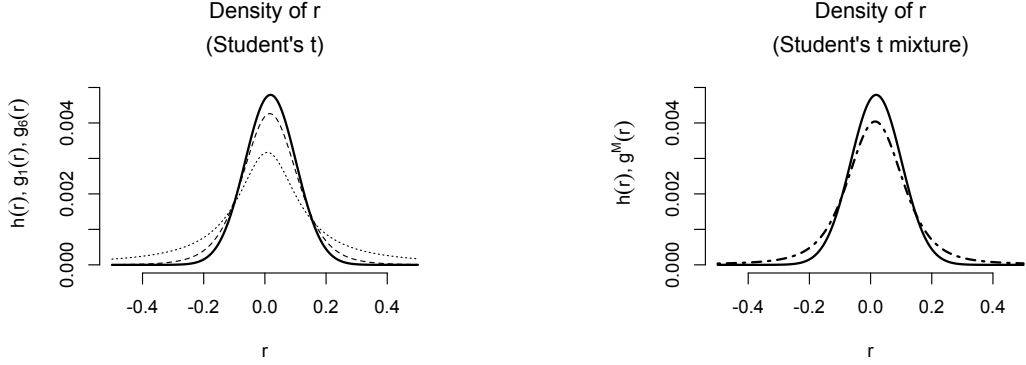


Figure 6: Empirical (solid) vs. posterior predictive (dotted and/or dashed) densities (parameterizations used are $\mathcal{N}(\hat{\mu}_r, \hat{\Sigma}_r)$ with $\hat{\mu}_r := 0.0184$ and $\hat{\Sigma}_r := 0.0069$ vs. $\mathcal{T}_1(\bar{r}_t, 2s_t^2)$ with $\bar{r}_t := 0.0082$ and $s_t^2 := 0.0210$ vs. $\mathcal{T}_{2\alpha_t}(\ell_t, \beta_t(1 + 1/\lambda_t)/\alpha_t)$ with $\ell_t := 0.0162$, $\lambda_t := 6$, $\alpha_t := 3$ and $\beta_t := 0.0630$ (left panel); $\mathcal{N}(\hat{\mu}_r, \hat{\Sigma}_r)$ as above vs. $g_r^M(\cdot)$ corresponding to $\pi^M(\cdot)$ from equation (20) with $p = 0.2$ as well as components $\mathcal{T}_1(\bar{r}_t, 2s_t^2)$ and $\mathcal{T}_{2\alpha_t}(\ell_t, \beta_t(1 + 1/\lambda_t)/\alpha_t)$ as above (right panel)).

this study implies the expected regime duration of 5 quarters,²⁷ we can plot the implied Student's t distributions $\mathcal{T}_6(\cdot)$ and $\mathcal{T}_1(\cdot)$, as well as the expected mixture distribution $\pi^M(\cdot)$ (see, respectively, left and right panels of Figure 6). The posterior mixture distribution “spreads out” in comparison to the observed sample, which is consistent with the pattern of probability weighting we saw in part §2 (Figure 2), as well as in §3.2 (Figures 4 and 5). Moreover, the result is effectively close to what we saw at the very outset in both panels of Figure 1: even though there is more probability mass in the central area of the distribution $\pi^M(\cdot)$ than in the distributions $g_r(\cdot)$ or $g_r^{KT}(\cdot)$, the remaining probability mass is allocated disproportionately more towards the far tails in the former than in the latter due to the presence of the Cauchy component in $\pi^M(\cdot)$.²⁸

Lastly, we note the following regularity.

Proposition 2 (“Complacency”). *With the progression of time t , where $t^* \leq t < t^{**}$, and an increasing distance from the last parameter change, the mixture distribution $\pi^M(\cdot)$ defined in equation (20) exhibits a compression of its tails, thus approaching the Normal distribution (although not converging, in distributional sense, to the Normal even in infinite time as long as $p > 0$).*

²⁷Assuming the length of a business cycle is about 5 years (in accordance with the NBER, see <http://www.nber.org/cycles.html>), each phase—such as trough, expansion, peak, contraction—lasts on average 5 quarters.

²⁸Precise calibration and matching in terms of the free parameter, change probability p , is not the purpose of our study and deserves a separate investigation.

Proof. A growing value of t increases, by equation (13), the magnitude of α_t , which in turn raises the number of degrees of freedom in the non-Cauchy component $\mathcal{T}_{2\alpha_t}(\ell_t, \beta_t(1 + 1/\lambda_t)/\alpha_t)$ of the Student’s t mixture from equation (20). □

As more observations from the new regime accumulate, the parameter posteriors sharpen and reduce the corresponding uncertainty (i.e., the dispersion and tail-thickness of the non-Cauchy component above), inducing a more confident behavior of the risk-averse investors. However, the chance of a regime change (and a switch to a Cauchy component) is always present. Once it happens, the parameter uncertainty will blow up (rapidly increasing the dispersion and tail-thickness of the mixture distribution), which to an outside observer (an econometrician) will look like a period of “complacency” followed by something reminiscent of “Minsky moment”²⁹ and an entailed crisis/realignment. That is, investors are aware of the chance of a regime change and (optimally) account for it, but they can not possibly predict and fully prepare for it.³⁰

To sum up, we have a highly stylized economic model that allows for a certain evolution over time: a Gaussian setup with underlying conditioning parameters that are stochastic, which is a primitive example of a non-ergodic stochastic environment.³¹ (Most of the time the random variables’ realizations are close to their previous values; but, on rare occasions, they deviate far from their recent past, potentially hitting the investors hard.)

Investors rationally respond to such parameter/state uncertainty, taking a Bayesian approach with the priors constructed by the methods of invariance and maximum entropy: an Inverse-Gamma prior for the variance, and a conditionally Normal prior for the mean. (This is a disciplined way of dealing with parameter uncertainty, and by construction accounts for undersampled or even unobserved events such as “rare disasters”.)

As a result, we obtain the posterior distributions of a Student’s t form (sometimes even its most extremal, Cauchy, version), although the true driving probabilities are Normal, conditionally on stochastic mean and variance. Empirically, the indications of a stochastic underlying parametric structure may be the long-known heavy tails/overdispersion (Mandelbrot, 1963; Fama, 1963) and state- or time-varying volatility (Officer, 1973; Black, 1976; Latané and Rendleman, 1976) of the asset returns. In turn, the low–degrees-of-freedom Student’s t distributions and mixtures thereof can help explain the asset pricing puzzles and regularities from parts §2.1 and §2.3.

Lastly, recognition of non-ergodicity immediately obviates the restriction (usually associated with a naïve form of “rational expectations”) that sample statistics, for a sufficiently long time-series, must converge to, or at least be informative of population statistics.

²⁹No relation to the mixture moments from Proposition 1. See, e.g., McCulley (2009).

³⁰Of course, a caveat applies that a constant change probability p is a very crude modeling shortcut, which ignores interesting fluctuations such as, for instance, endogenous dynamics of financial imbalances.

³¹Existing works recognizing the limitations of the available sample of macro-financial statistics and/or the evolving nature of the underlying economic parameters are the mentioned above Weitzman (2007), but also McGrattan and Prescott (2003 and 2005), Dimson et al. (2003), Fama and French (2002).

Thus, the mixture distribution $\pi^M(\cdot)$ above never converges to the population distribution of returns r_t : conditionally on the current regime parameters, it does not because of the overarching presence of the risk of regime change; while unconditionally, it does not because the population itself is in some sense an ill-defined object in these circumstances.

4 Discussion

The theoretical arguments in part §3 imply that the prospect theory’s probability weighting transformations and the large-variance Normal distribution we have seen earlier in §2 may be just approximations of the (means of) Beta and (a mixture of) the Student’s t posterior densities that arise from a conventional Bayesian updating procedure under certain (discipline about) prior assumptions. In other words, these theoretical results suggest a **“rationalization” of the stylized empirical facts** such as Allais’ paradox and Kahneman-Tversky’s probability distortions as well as the equity premium and risk-free rate puzzles.

We use the same conceptual approach for modeling both the lottery choices in the lab experiments as well as the asset prices in financial markets. Nevertheless, given the obvious differences in problem setups, strictly speaking there would be little reason for the results in these two domains to match not just qualitatively, but also quantitatively (i.e., for the same amount of probability reweighting as observed in the lottery experiments to work for asset prices as well). Note, however, that in the application to binary lotteries, we use a very elementary Bayesian inference approach; while in the application to asset prices, we use a fully developed general equilibrium model, which in particular, allows for an extra free parameter, the coefficient of relative risk aversion γ .

Although there are **alternative models** aimed at rationalizing the above empirical facts, they are not applicable to both micro-level (probability distortions, e.g., in laboratory experiments with lottery choices) and macro-level phenomena (risk premia in asset prices on financial markets).

Micro-level theories, such as various axiomatic derivations of rank-dependent utility (an early overview is available in Wakker, 1994), those based on the costs of self-control (Fudenberg and Levine, 2006; 2011) or noise in neural information processing (Steiner and Stewart, 2016) rely on specific constraints stemming from individual preferences and costs, which in competitive and liquid markets, presumably should be “arbitraged away” by professional/institutional—i.e., risk-neutral, relatively immune to idiosyncratic errors, etc.—agents.

On the other hand, macro-level theories emphasizing the role of uncertainty about economy’s evolving structure (Weitzman, 2007), habit formation (Campbell and Cochrane, 1999) and so on do not directly translate to static stochastic choice setups. (See Appendix §A for more details.)

Next, the recognition that learning and inference from the limited available informa-

tion are inherent constituents of decision-making under risk offers a **rational unifying framework** for a theoretical normative axiomatic approach to and empirical positive description of the optimal choice under risk. Indeed, rational behavior of the decision-makers entails that their optimal decisions, made effectively under the unobserved but inferred population distribution, obviously satisfy the standard von Neumann-Morgenstern axioms; however, viewed by experimenters/econometricians under an observed but limited sample distribution, these same decisions seem to deviate from rationality and exhibit phenomena such as the Allais paradox or the equity-premium puzzle.

Lastly, here we obtain a manifestation of the well-known **shrinkage** phenomenon. Originally it was proposed for estimation of the mean of a multivariate Normal distribution from an ultra-small sample of a single vector observation (Stein, 1956; James and Stein, 1961).³² The concept of Stein-type shrinkage stimulated the reinterpretation of old as well as the development of new methods of statistical analysis such as the Good–Turing frequency estimation in computational linguistics and machine learning (Good, 1953); least absolute shrinkage and selection operator, or Lasso, estimation in statistics and machine learning (Tibshirani, 1996); improved estimates of the variance-covariance matrix (Ledoit and Wolf, 2004; Jagannathan and Ma, 2003) and of the mean vector (Jorion, 1985; 1986) of portfolio returns in finance.

In our case it amounts to deforming (“shrinking”) the distribution of interest in the direction towards the invariant ignorance prior, i.e., a quasi-uniform distribution (“shrinkage target”). Indeed, this quasi-uniform prior distribution serves as an “Occam factor” that favors simple models. In turn, such an adjustment reduces the chances of overfitting the observed data, which would have been undesirable in the situation of undersampling. Specifically, distributions of the lotteries’ (transformed) probabilities are shrunk towards the uniform distribution on the $[0, \pi/2]$ interval, thus biasing the small probabilities upward and the large probabilities downward. Distributions of the asset returns’ (log-) variances and (conditional) means are both shrunk towards the uniform distributions on \mathbb{R} , prudently³³ amplifying the uncertainty about these returns-driving parameters. (See Appendix §B for more details.)

5 Conclusion

This paper documents some stylized facts about the paradoxical choice behavior humans exhibit in laboratory experiments and about the asset pricing puzzles in financial markets, as well as presents the empirical connection between these two decision-making domains. Intuitively, the above stylized facts can be understood as being motivated by the need to account for “rare” unobserved contingencies when working with only “small” samples of the relevant data (e.g., when facing new gambles or new economic regimes).

³²In this connection, also see regularization, e.g., in Chen and Haykin (2002), Bickel and Li (2006).

³³For Good reason, if you will.

We can rationalize such decision-making behavior by taking a dynamic perspective and appealing to Bayesian updating under the requirements of parameterization invariance and maximum entropy. Effectively, this formulates rational optimality-motivated foundations behind the probability distortions stipulated by the (cumulative) prospect theory of Kahneman and Tversky in microeconomics and decision theory as well as behind the equity premium/risk-free rate and some other puzzles in macro-finance.

References

- [1] Abdellaoui, Mohammed. (2000) “Parameter-Free Elicitation of Utility and Probability Weighting Functions”, *Management Science*, 46(11): 1497–1512.
- [2] Ang, Andrew and Allan Timmermann. (2012) “Regime Changes and Financial Markets”, *Annual Review of Financial Economics*, 4: 313–337.
- [3] Allais, Maurice. (1953) “Le Comportement de l’Homme Rationnel devant le Risque: Critique des Postulats et Axiomes de l’Ecole Americaine”, *Econometrica*, 21(4): 503–546.
- [4] Bansal, Ravi and Amir Yaron. (2004) “Risks for the Long-Run: A Potential Resolution of Asset Pricing Puzzles.” *Journal of Finance*, 59(4), 1481–1509.
- [5] Barberis, Nicholas and Ming Huang. (2008) “Stocks as Lotteries: The Implications of Probability Weighting for Security Prices”, *American Economic Review*, 98(5): 2066–2100.
- [6] Barro, Robert J. (2006) “Rare Disasters and Asset Markets in the Twentieth Century.” *The Quarterly Journal of Economics*, 121(3): 823–866.
- [7] Berg, Christian, and Christophe Vignat. (2010) “On the Density of the Sum of Two Independent Student t-random Vectors”, *Statistics and Probability Letters*, 80: 1043–1055.
- [8] Bickel, Peter J. and Bo Li (2006) “Regularization in Statistics.” *Test*, 15(2): 271–344.
- [9] Black, F. (1976) “Studies of Stock Price Volatility Changes”, Proceedings of the 1976 Meeting of the Business and Economic Statistics Section, American Statistical Association, Washington DC: 177–181.
- [10] Blume, Marshall E. and Irwin Friend (1975) “The Asset Structure of Individual Portfolios and Some Implications for Utility Functions”, *The Journal of Finance*, 30(2): 585–603.
- [11] Bossaerts, Peter, Kerstin Preuschoff and Ming Hsu. (2008). “The Neurobiological Foundations of Valuation in Human Decision Making Under Uncertainty.” In: Paul W. Glimcher, Colin F. Camerer, Ernst Fehr, Russell A. Poldrack (eds.), *Neuroeconomics: Decision Making and the Brain*, Academic Press.
- [12] Bruhin, Adrian, Helga Fehr-Duda and Thomas Epper. (2010) “Risk and Rationality: Uncovering Heterogeneity in Probability Distortion”, *Econometrica*, 78(4): 1375–1412.
- [13] Camerer, Colin F. and Teck-Hua Ho. (1994) “Violations of the Betweenness Axiom and Nonlinearity in Probability”, *Journal of Risk and Uncertainty*, 8(2): 167–196.
- [14] Campbell, John Y. and John Cochrane. (1999) “By Force of Habit: A Consumption-Based Explanation of Aggregate Stock Market Behavior.” *Journal of Political Economy*, 107: 205–251.
- [15] Chen, Zhe and Simon Haykin. (2002). “On Different Facets of Regularization Theory.” *Neural Computation*, 14: 2791–2846.
- [16] Chernoff, Herman and Lincoln E. Moses. (1959) *Elementary Decision Theory*, New York: Wiley.
- [17] DeGiorgi, Enrico G. and Shane Legg. (2012) “Dynamic portfolio choice and asset pricing with narrow framing and probability weighting”, *Journal of Economic Dynamics and Control*, 36: 951–972.
- [18] DeGroot, Morris H. (1970) *Optimal Statistical Decisions*, McGraw-Hill.
- [19] Dimson, Elroy, Paul Marsh and Mike Staunton. (2003) “Global Evidence on the Equity Risk Premium”. *LBS Institute of Finance and Accounting Working Paper*, No. IFA 385.
- [20] Epstein, Larry G. and Stanley E. Zin. (1989) “Substitution, Risk Aversion, and the Temporal Behavior of Consumption and Asset Returns: A Theoretical Framework.” *Econometrica*, 57: 937–969.

- [21] Epstein, Larry G. and Stanley E. Zin. (1990) “‘First-Order’ Risk Aversion and the Equity Premium Puzzle.” *Journal of Monetary Economics*, 26(3): 387–407.
- [22] Epstein, Larry G. and Stanley E. Zin. (1991) “Substitution, Risk Aversion, and the Temporal Behavior of Consumption and Asset Returns: An Empirical Analysis.” *Journal of Political Economy*, 99(2): 263–286.
- [23] Fama, Eugene F. (1963) “Mandelbrot and the Stable Paretian Hypothesis.” *The Journal of Business*, 36(4): 420–429.
- [24] Fama, Eugene F. and Kenneth R. French. (2002) “The Equity Premium”, *The Journal of Finance*, 57: 637–659.
- [25] Fehr-Duda, Helga and Thomas Epper. (2012) “Probability and Risk: Foundations and Economic Implications of Probability-Dependent Risk Preferences”, *Annual Review of Economics*, 4: 567–593.
- [26] Fox, Craig R., Brett A. Rogers and Amos Tversky. (1996) “Options Traders Exhibit Subadditive Decision Weights”, *Journal of Risk and Uncertainty*, 13(1): 5–17.
- [27] French, Kenneth R. and James M. Poterba. (1991) “Investor Diversification and International Equity Markets”, *American Economic Review*, 81(2): 222–226.
- [28] Fudenberg, Drew, and David K. Levine. (2006). “A Dual-Self Model of Impulse Control”, *American Economic Review*, 96(5), 1449–1476.
- [29] Fudenberg, Drew, and David K. Levine. (2011). “Risk, Delay, and Convex Self-Control Costs”, *American Economic Journal: Microeconomics*, 3(3), 34–68.
- [30] Gelman, Andrew, John B. Carlin, Hal S. Stern, Donald B. Rubin. (2004) *Bayesian data analysis*, 2nd edition, Chapman & Hall.
- [31] Good, I.J. (1953). “The Population Frequencies of Species and the Estimation of Population Parameters”, *Biometrika*, 40(3/4): 237–264.
- [32] Hamilton, James D. (2016) “Macroeconomic Regimes and Regime Shifts.” In: John B. Taylor and Harald Uhlig (eds.), *Handbook of Macroeconomics*, volume 2, part A, chapter 3, pages 163–201, Elsevier.
- [33] Hansen, Lars P. (2007) “Beliefs, Doubts and Learning: Valuing Macroeconomic Risk.” *American Economic Review*, 97 (2): 1–30.
- [34] Hansen, Lars P. (2014) “Nobel Lecture: Uncertainty Outside and Inside Economic Models”, *Journal of Political Economy*, 122(5): 945–987.
- [35] Hansen, Lars P., John C. Heaton and Nan Li. (2008) “Consumption Strikes Back?: Measuring Long-Run Risk.” *Journal of Political Economy*, 116(2): 260–302.
- [36] Hansen, Lars Peter and Ravi Jagannathan. (1991) “Implications of Security Market Data for Models of Dynamic Economies”, *Journal of Political Economy*, 99(2): 225–262.
- [37] Hansen, Lars P. and Thomas J. Sargent. (2007b) “Recursive Robust Estimation and Control Without Commitment.” *Journal of Economic Theory*, 136 (1): 1–27.
- [38] Hausser, Jean and Korbinian Strimmer. (2009) “Entropy Inference and the James-Stein Estimator, with Application to Nonlinear Gene Association Networks”, *Journal of Machine Learning Research*, 10: 1469–1484.
- [39] Horn, Roger A. and Charles R. Johnson. (1985) *Matrix Analysis*, Cambridge University Press.
- [40] Jagannathan, Ravi and Tongshu Ma. (2003) “Risk Reduction in Large Portfolios: Why Imposing the Wrong Constraints Helps.” *Journal of Finance*, 58(4): 1651–1684.

- [41] James, W. and Charles M. Stein. (1961) “Estimation with Quadratic Loss”, Proc. Fourth Berkeley Symp. Math. Statist. Prob. 1: 361–379.
- [42] Jaynes, Edwin Thompson. (1957a) “Information Theory and Statistical Mechanics”, *Physical Review*, 106(4): 620–630.
- [43] Jaynes, Edwin Thompson. (1957b) “Information Theory and Statistical Mechanics. II”, *Physical Review*, 108(2): 171–190.”
- [44] Jaynes, Edwin T. (2003) *Probability Theory: The Logic of Science*, Cambridge University Press.
- [45] Jeffreys, Harold. (1939) *Theory of Probability*, Oxford University Press.
- [46] Jeffreys, Harold. (1946) “An Invariant Form for the Prior Probability in Estimation Problems,” *Proceedings of the Royal Society A*, 186(1007): 453–461.
- [47] Johnstone, Ian M. (2001) “On the Distribution of the Largest Eigenvalue in Principal Components Analysis.” *Annals of Statistics*, 29(2): 295–327.
- [48] Jorion, Philippe. (1985) “International Portfolio Diversification with Estimation Risk.” *Journal of Business*, 58(3): 259–278.
- [49] Jorion, Philippe. (1986) “Bayes-Stein Estimation for Portfolio Analysis.” *Journal of Financial and Quantitative Analysis*, 21(3): 279–292.
- [50] Kahneman, Daniel and Amos Tversky. (1979) “Prospect Theory: An Analysis of Decision under Risk”, *Econometrica*, 47(2): 263–291.
- [51] Kass, Robert E. and Larry Wassermann. (1996) “The Selection of Prior Distributions by Formal Rules”, *Journal of the American Statistical Association*, 91(435): 1343–1370.
- [52] Kelly, Morgan. (1995) “All Their Eggs in One Basket: Portfolio Diversification of US Households”, *Journal of Economic Behavior and Organization*, 27: 87–96.
- [53] Keynes, John Maynard. (1921) *A Treatise on Probability*, London: Macmillan & Co.
- [54] Kliger, Doron and Ori Levy. (2009) “Theories of Choice Under Risk: Insights from Financial Markets”, *Journal of Economic Behavior and Organization*, 71: 330–346.
- [55] Laplace, Pierre-Simon. (1825) *Philosophical Essay on Probabilities*, London: Chapman & Hall.
- [56] Latané, Henry A. and Richard J. Rendleman Jr. (1976) “Standard Deviations of Stock Price Ratios Implied in Option Prices”, *Journal of Finance*, 31(2): 369–381.
- [57] Ledoit, Olivier and Michael Wolf (2004) “A Well-Conditioned Estimator for Large-Dimensional Covariance Matrices.” *Journal of Multivariate Analysis*, 88(2): 365–411.
- [58] Ledoit, Olivier and Michael Wolf. (2013) “Spectrum Estimation: A Unified Framework for Covariance Matrix Estimation and PCA in Large Dimensions.” Manuscript.
- [59] List, John A. (2004) “Neoclassical Theory Versus Prospect Theory: Evidence from the Marketplace”, *Econometrica*, 72(2): 615–625.
- [60] List, John A. and Michael S. Haigh. (2005) “A Simple Test of Expected Utility Theory Using Professional Traders”, *Proceedings of the National Academy of Sciences*, 102(3): 945–948.
- [61] Lucas, Robert E., Jr. (1978) “Asset Prices in an Exchange Economy.” *Econometrica*, 46 (6): 1429–1445.
- [62] MacKay, David J. C. (2003) *Information Theory, Inference, and Learning Algorithms*, Cambridge University Press.
- [63] Malloy, Christopher J., Tobias J. Moskowitz and Annette Vissing-Jørgensen (2009). “Long-Run Stockholder Consumption Risk and Asset Returns.” *Journal of Finance*, 64(6): 2427–2479.

- [64] Mandelbrot, Benoit. (1963) “The Variation of Certain Speculative Prices.” *The Journal of Business*, 36(4): 394–419.
- [65] McCulley, Paul. (2009) “The Shadow Banking System and Hyman Minsky’s Economic Journey”, *PIMCO Global Central Bank Focus*, May.
- [66] McGrattan, Ellen R. and Edward C. Prescott (2003) “Average Debt and Equity Returns: Puzzling?” *American Economic Review*, 93: 392–397.
- [67] McGrattan, Ellen R. and Edward C. Prescott (2005) “Taxes, Regulations, and the Value of U.S. and U.K. Corporations”, *Review of Economic Studies*, 92: 767–796.
- [68] Mehra, Rajnish and Edward C. Prescott. (1985) “The Equity Premium: A Puzzle”, *Journal of Monetary Economics*, 15(2): 145–161.
- [69] Menzly, Lior, Tano Santos and Pietro Veronesi. (2004) “Understanding Predictability.” *Journal of Political Economy*, 112(1): 1–47.
- [70] Minka, Thomas P. (2001) “Inferring a Gaussian distribution.” Manuscript.
- [71] Murphy, Kevin P. (2012) *Machine Learning: A Probabilistic Perspective*, The MIT Press
- [72] Nemenman, Ilya, Fariel Shafee and William Bialek. (2002) “Entropy and Inference, Revisited,” In: T. G. Dietterich, S. Becker, and Z. Ghahramani (eds.), *Advances in Neural Information Processing Systems*, 14, MIT Press.
- [73] Officer, R. R. (1973) “The Variability of the Market Factor of the New York Stock Exchange.” *The Journal of Business*, 46(3): 434–453.
- [74] Orlitsky, Alon, Narayana P. Santhanam, Junan Zhang. (2003) “Always Good Turing: Asymptotically Optimal Probability Estimation”, *Science*, 302(5644): 427–431.
- [75] Paninski, Liam. (2003) “Estimation of Entropy and Mutual Information”, *Neural Computation*, 15: 1191–1253.
- [76] Parker, Jonathan A. and Christian Julliard. (2005) “Consumption Risk and the Cross Section of Expected Returns.” *Journal of Political Economy*, 113(1): 185–222.
- [77] Polkovnichenko, Valery. (2005) “Household Portfolio Diversification: A Case for Rank-Dependent Preferences”, *Review of Financial Studies*, 18(4): 1467–1502.
- [78] Polkovnichenko, Valery and Feng Zhao. (2013) “Probability Weighting Functions Implied in Options Prices”, *Journal of Financial Economics*, 107: 580–609.
- [79] Rietz, Thomas. (1988) “The Equity Risk Premium: A Solution.” *Journal of Monetary Economics*, 22: 117–131.
- [80] Rissanen, Jorma. (1983) “A Universal Prior for Integers and Estimation by Minimum Description Length”, *The Annals of Statistics*, 11(2): 416–431.
- [81] Routledge, Bryan R. and Stanley E. Zin. (2010) “Generalized Disappointment Aversion and Asset Prices”, *Journal of Finance*, 65(4): 1303–1332.
- [82] Statman, Meir. (1987) “How Many Stocks Make a Diversified Portfolio?”, *The Journal of Financial and Quantitative Analysis*, 22(3): 353–363.
- [83] Stein, Charles M. (1956), “Inadmissibility of the Usual Estimator for the Mean of a Multivariate Distribution”, Proc. Third Berkeley Symp. Math. Statist. Prob. 1: 197–206.
- [84] Stein, Charles M. (1975) “Estimation of a Covariance Matrix”, Rietz Lecture, 39th Annual Meeting of the Institute of Mathematical Statistics. Atlanta, GA.

- [85] Stein, Charles M. (1986) “Lectures on the Theory of Estimation of Many Parameters.” *Journal of Soviet Mathematics*, 34(1): 1373–1403.
- [86] Steiner, Jakub and Colin Stewart (2016) “Perceiving Prospects Properly,” *American Economic Review*, 106(7): 1601–1631.
- [87] Tibshirani, Robert. (1996) “Regression Shrinkage and Selection via the Lasso”, *Journal of the Royal Statistical Society. Series B*, 58(1): 267–288.
- [88] Tsai, Jerry and Jessica A. Wachter. (2016) “Rare Booms and Disasters in a Multisector Endowment Economy”, *Review of Financial Studies*, 29(5): 1113–1169.
- [89] Tversky, Amos and Daniel Kahneman. (1992) “Advances in Prospect Theory: Cumulative Representation of Uncertainty”, *Journal of Risk and Uncertainty*, 5: 297–323.
- [90] Valiant, Gregory and Paul Valiant. (2017) “Estimating the Unseen: Improved Estimators for Entropy and other Properties”, *Journal of the ACM*, 64(6)((37)): 1–41.
- [91] Veronesi, Pietro. (2004) “The Peso Problem Hypothesis and Stock Market Returns”, *Journal of Economic Dynamics and Control*, 28: 707–725.
- [92] Wakker, Peter. (1994) “Separating Marginal Utility and Probabilistic Risk Aversion”, *Theory and Decision*, 36: 1–44.
- [93] Weil, Philippe. (1989) “The Equity Premium Puzzle and the Risk-Free Rate Puzzle.” *Journal of Monetary Economics*, 24: 401–421.
- [94] Weitzman, Martin L. (2007) “Subjective Expectations and Asset-Return Puzzles.” *American Economic Review*, 97 (4): 1102–1130.
- [95] Wu, George and Richard Gonzalez. (1996) “Curvature of the Probability Weighting Function”, *Management Science*, 42(12): 1676–1690.
- [96] Ziegler, Alexandre. (2007) “Why Does Implied Risk Aversion Smile?” *The Review of Financial Studies*, 20(3): 859–904.