

Avoiding Risk and Avoiding Evidence

Catrin Campbell-Moore and Bernhard Salow

February 28, 2017

Abstract

It is natural to think that there's something epistemically objectionable about avoiding evidence, at least in ideal cases. We argue that this natural thought is inconsistent with a kind of risk avoidance that is both wide-spread and intuitively rational. More specifically, we argue that if the kind of risk avoidance defended by Buchak (2013) is rational, avoiding evidence can be epistemically commendable.

In the course of our argument we also lay some foundations for studying epistemic utility, or accuracy, when considering risk-avoidant agents.

Is it ever reasonable not to gather available evidence? Sure it is. Gathering and processing the evidence is almost never free, costing you time and cognitive energy if nothing else. Even if it were free, you might know that the evidence doesn't bear on anything important. And even if the evidence were both free and relevant, you might still be worried that you'll misevaluate it.

But what about cases in which none of these worries arise? *Ideal cases*, in which (i) gathering the evidence incurs no cost whatsoever, (ii) the evidence is potentially relevant, and (iii) you're certain to process it rationally?¹ In these cases, it seems wrong to ignore evidence. Classical decision theory agrees: it says that not gathering the evidence is *instrumentally irrational*, a bad way of pursuing your goals.² However, decision theories that allow for a kind of 'risk avoidance' not permitted by classical decision theory do not.³ Classical and risk-avoidant decision theories thus disagree about:

¹Evidence is potentially relevant if the agent has positive credence that it makes a difference to what outcome she achieves; in the instrumental case, this means that the agent has positive credence that the evidence will make her choose a different option in some decision problem she will face, while in the epistemic case it means that the agent has positive credence that it will make a difference to what credence she assigns to some claim that matters.

One might worry that, because gathering evidence always carries *some* cost, there couldn't be any ideal cases, making our principles vacuous. Perhaps that is so. However, given a precise definition of costs, one could define what it is for the costs of gathering the evidence to be negligible, compared to the other things at stake – one could then redefine ideal cases as ones where the costs are negligible. However, we can't give a precise definition of 'costs' (and hence of cost-negligibility) here, and so will opt for the simpler formulation. We should, however, note that, when considering epistemic principles, the relevant costs should be understood as *epistemic* ones – for example, the loss of time, energy, or opportunity to gather other evidence.

²Good (1967)

³Wakker (1988)

Look-I. In ideal cases, one is *instrumentally* required to gather the evidence.

It thus looks as though evidence-avoidance can be rational (in ideal cases) if and only if risk-avoidance is.

Recently, however, Buchak (2010) has pointed out that this reasoning ignores an important distinction. For, even if avoiding evidence is rational from a purely instrumental perspective, it might nonetheless be irrational from an epistemic one. This would mean that instrumental and epistemic rationality sometimes conflict. But that is independently plausible: if you offer to pay me a lot of money to take a pill that will instil in me the (mistaken) belief that I had eggs for breakfast today, I may be instrumentally required to accept but epistemically required to decline. It thus seems as though, in rejecting Look-I, advocates of risk-avoidance might nonetheless be able to accept:

Look-E. In ideal cases, one is *epistemically* required to gather the evidence.

If that's right, then even if avoiding risk is rational, there remains an important sense in which avoiding evidence is not.

In this paper, we close this gap by arguing that if risk-avoidance is rational, then avoiding evidence is sometimes epistemically rational (even in ideal cases). More precisely, we argue that if the kind of risk-avoidance permitted by Buchak's decision theory is rational then it can be epistemically rational to avoid evidence even in ideal cases. We leave open whether the correct reaction is to reject the rationality of risk aversion, or to conclude that, even from an epistemic point of view, there needn't be anything wrong with avoiding evidence.⁴

The plan is as follows. Sections 1 and 2 present relevant background, explaining the connection between risk-avoidance and Look-I, as well as Buchak's risk-sensitive decision theory. Section 3 begins the main argument, by making the case that epistemic rationality requires an agent to gather (or avoid) evidence if doing so is conducive towards securing epistemic goods; and hence that risk-sensitive agents are required to gather (or avoid) evidence if doing so maximizes risk-weighted expected epistemic utility. Section 4 discusses how to measure epistemic utility when working with risk-sensitive agents. Section 5 presents the core result showing that Look-E fails for risk avoidant agents. Sections 6 and 7 show that much of this result goes through even if certain assumptions (about how epistemic utility is measured, and how rational agents revise their beliefs) are significantly weakened. Section 8 sums up.

1 Risk aversion and Look-I

Classical decision theory requires agents to maximize expected utility. This means that all rational aversion to risky options must be reflected in the agent's

⁴Strictly speaking, another option is to maintain that some risk aversion is rational, but not the kind described by Buchak's theory. We will offer some (non-*ad hominem*) reasons to focus on Buchak's theory in evaluating the interaction between risk avoidance and Look-E; but, those reasons won't be conclusive, and so this third option remains open.

utility function: if it's rational to prefer a sure \$5 over a bet that pays \$10 if a fair coin lands heads (and nothing otherwise), this is only because money has diminishing marginal utility, so that the utility of winning \$10 is less than twice that of securing \$5. But many have argued that this does not vindicate all the ways in which rational agents can be risk averse. Perhaps the best-known example is that adjusting the utilities doesn't allow us to capture the 'Allais Preferences'.⁵ Consider the following four lotteries

	Ticket 1	Ticket 2-11	Ticket 12-100
L_1	\$1,000	\$1,000	\$1,000
L_2	\$0	\$2,000	\$1,000
L_3	\$1,000	\$1,000	\$0
L_4	\$0	\$2,000	\$0

Note that L_1 differs from L_3 in exactly the same way as L_2 differs from L_4 ; in both cases, we've just replaced 88 \$1,000-tickets with losing tickets. This means that expected utility theory predicts that, regardless of how you value money, if you prefer L_1 to L_2 , you also prefer L_3 to L_4 . However, many people report preferences contradicting this, favouring L_1 over L_2 but L_4 over L_3 . Moreover, such preferences seem to make sense. L_1 feels preferable to L_2 in part because it is risk-free, while L_3 has no such attraction over L_4 – that L_1 differs from L_3 in the same way as L_2 differs from L_4 just hides this intuitively relevant point.

Cases like these have motivated the development of decision theories that allow for other forms of risk-avoidance, and thus permit the Allais preferences. Rejecting classical decision theory, however, means that we need to reexamine its treatment of evidence-gathering. I.J. Good (1967) famously proved that, when faced with a choice between either (i) choosing an option now or (ii) gathering some cost-free evidence, conditionalizing on it, and then choosing the option that maximizes expected utility relative to the updated credences, (ii) will have a strictly higher expected utility than (i) whenever the new evidence might lead one to choose a different option. If we follow classical decision theory in assuming that instrumental rationality requires us to maximize expected utility, this establishes Look-I. But if we reject classical decision theory, it obviously does not.

Moreover, no theory which permits the type of risk avoidance manifest in the Allais preferences can recover an analogous result.⁶ We won't go over the argument here. But, very roughly, the idea is that, since the choice between L_1 and L_2 on the one hand, and between L_3 and L_4 on the other, is the same if you know that your ticket is 1–11 (and is inconsequential if you know that your ticket is 12–100), you are bound to 'switch' preferences in one of the two cases upon finding out whether your ticket is 1–11 or 12–11. But this means

⁵Allais (1953). The preferences usually involve \$1m where we have \$1,000 and \$5m where we have \$2,000. These changes affect nothing of substance; they just simplify the calculations in section 2.

⁶Cf Buchak (2013, p.171-173). Wakker (1988) establishes the even stronger claim that no theory which allows for any violations of the 'Sure-Thing Principle' (a principle of classical decision theory violated by the Allais preferences) can recover Look-I.

that, in that case, you will pick an option you (initially) consider sub-optimal if you gather the evidence before choosing.

This result is somewhat strange, but Buchak offers a helpful diagnosis. We can call some evidence *instrumentally misleading* if it makes it rational to perform an action that is, as a matter of fact, worse than the one you would have performed if you hadn't received that evidence. Since learning that your ticket is 1–11 makes it rationally permissible to perform a different action (taking L_2 instead of L_1), it has a risk of being misleading; and since there is a $1/11$ chance of your ticket being ticket 1 (and hence L_1 still having the better outcome) even if it is one of 1–11, this risk is substantial. Of course, that risk must be weighed against the fact that learning 1–11 could be *instrumentally 'truth-guiding'*, making it rational to perform an action which in fact leads to a better outcome than you would otherwise achieve (taking L_2 instead of L_1 when you hold ticket 2–11). For a risk-neutral agent, these benefits always outweigh the risk. But for a risk-sensitive agent, they may not.

All theories which depart from expected utility theory to allow for intuitively rational forms of risk avoidance of the kind represented by the Allais Preferences thus reject Look-I. However, as Buchak (2010) points out, this leaves open that we may nonetheless always be *epistemically* required to gather evidence. Our aim in this paper is to investigate whether, if risk avoidance is rational, we are so required. To do that, we will need to employ a particular theory of risk-avoidance; for reasons we will explain shortly, we will focus on the theory developed by Buchak herself.

2 Buchak's Theory

On Buchak's theory, an agent is represented as having not just credences c and utilities U , but also a risk function r capturing her attitude towards risk – specifically how much she engages in worst-case-scenario style reasoning.

Mathematically, r is just an increasing function from $[0,1]$ to $[0,1]$, such that $r(0) = 0$ and $r(1) = 1$. To understand its role in calculating risk-weighted expected utilities note that ordinary risk-neutral expected utilities can be written, slightly non-standardly, as follows:

Risk neutral expected utility. Suppose an act, A , leads to outcomes o_1, \dots, o_n in states s_1, \dots, s_n , with $U(o_1) \leq \dots \leq U(o_n)$. Then

$$\begin{aligned} \text{Exp}_c U(A) &= U(o_1) \\ &\quad + (c(s_2) + \dots + c(s_n)) \cdot (U(o_2) - U(o_1)) \\ &\quad + \dots \\ &\quad + c(s_n) \cdot (U(o_n) - U(o_{n-1})) \end{aligned}$$

Intuitively, the (risk-neutral) expected utility is here calculated by first taking the utility of the worst-case scenario; adding the improvement over the worst-case scenario secured in the second-worst-case scenario, weighted by the

probability of securing at least that improvement; adding the improvement over the second-worst-case scenario secured the third-worst-case scenario, weighted by the probability of securing at least *that* improvement; and continuing like this until all the possible improvements have been taken into account. The risk-weighted expected utilities are calculated in exactly the same way, except that r can modify the weight which possible improvements receive:

Risk weighted expected utility. Suppose an act, A , leads to outcomes o_1, \dots, o_n in states s_1, \dots, s_n with $U(o_1) \leq \dots \leq U(o_n)$. Then:

$$\begin{aligned} \text{RExp}_c^r U(A) &= U(o_1) \\ &\quad + r(c(s_2) + c(s_3) + \dots + c(s_n)) \cdot (U(o_2) - U(o_1)) \\ &\quad + r(c(s_3) + \dots + c(s_n)) \cdot (U(o_3) - U(o_2)) \\ &\quad + \dots \\ &\quad + r(c(s_n)) \cdot (U(o_n) - U(o_{n-1})) \end{aligned}$$

If an agent has the risk function $r(x) = x$, this is equivalent to the usual formula.

But if r is more interesting, we can get different results. Suppose, for example, that $r(x) < x$. Then the weight given to the potential improvements will be less than it is in the expected utility calculation; and so the relative weight given to the worst-case scenario is increased. This means, for example, that one can have utilities that are linear with money and still prefer the sure \$5 (which has an $\text{RExp}_c^r U$ of 5) over the possible \$10 (which has an $\text{RExp}_c^r U$ of $0 + r(0.5) \cdot (10 - 0) < 5$) – one thus prefers the safe choice not because one values the ‘second’ \$5 less, but because one gives more weight to the worst-case scenario when evaluating one’s options.

Intuitively, risk-aversion is a matter not just of giving extra weight to the *worst*-case scenario, but more generally a matter of giving more (relative) weight to *worse* scenarios. The property mentioned above, that $r(x) < x$, is not quite enough to ensure this; but a slightly stronger property – convexity – is.⁷ Following Buchak, $r(x) = x^2$ will be our main example of a risk-profile which has this feature.

⁷ r is convex if $r(\lambda x + (1 - \lambda)y) > \lambda r(x) + (1 - \lambda)r(y)$ for $x \neq y$ and $\lambda \in (0, 1)$. To explain why this is the required property, we first introduce the abbreviation ‘ p_i ’ for the probability that we achieve at least outcome o_i (so that $p_i = \sum_{i \leq j \leq n} c(s_j)$), with $p_{n+1} := 0$ as a limiting case. Then we can rearrange the formulas for EU and REU as

$$\begin{aligned} \text{RExp}_c^r U(A) &= \sum_{i=1 \dots n} U(s_i) \cdot (r(p_i) - r(p_{i+1})) \\ \text{Exp}_c U(A) &= \sum_{i=1 \dots n} U(s_i) \cdot (p_i - p_{i+1}) \end{aligned}$$

So we can think of the weight given to an outcome o_i in the REU calculation as the weight it receives in the EU calculation (namely $p_i - p_{i+1}$), but scaled by the factor

$$\frac{r(p_i) - r(p_{i+1})}{p_i - p_{i+1}}.$$

This factor can either increase or decrease the weight which this outcome has in determining the REU (relative to the weight it has in determining the EU). To say that worse outcomes receive additional weight is then to say that the scaling factors will become smaller as that

Such an agent is therefore inclined to avoid risks in a way not represented in her utilities.⁸

The presence of the risk profile also allows the theory to rationalize the Allais preferences described earlier. Recall the lotteries we mentioned, and consider an agent for whom utility and money are interchangeable, but who is risk-avoidant in line with the risk profile $r(x) = x^2$:

	Ticket 1	Ticket 2-11	Ticket 12-100
L_1	1,000	1,000	1,000
L_2	0	2,000	1,000
L_3	1,000	1,000	0
L_4	0	2,000	0

Since 1000 is L_1 's worst-case scenario, and there is no possibility of improvement, our agent assigns

$$\text{RExp}_c^r U(L_1) = 1000.$$

For L_2 , the $\text{RExp}_c^r U$ is given by looking at the base-line of 0 if the ticket is ticket 1; factoring in the possible improvement of securing 1000 more than that if the ticket is 12–100; and then factoring in the possible improvement of securing an additional 1000 if the ticket is 1–11. Plugging in the probabilities and risk-profile, this amounts to

$$\begin{aligned} \text{RExp}_c^r U(L_2) &= 0 + (0.99)^2 \cdot (1000 - 0) + (0.1)^2 \cdot (2000 - 1000) = 990.1. \\ &< \text{RExp}_c^r U(L_1) \end{aligned}$$

Similar reasoning shows:

$$\begin{aligned} \text{RExp}_c^r U(L_3) &= 0 + (0.11)^2 \cdot (1000 - 0) = 12.1 \\ &< \text{RExp}_c^r U(L_4) &= 0 + (0.1)^2 (2000 - 0) = 20 \end{aligned}$$

The agent thus has the preferences which the classical theory couldn't capture, preferring L_1 to L_2 but L_4 to L_3 .

Moreover, if our agent were to learn that her ticket was in the range 1–11, the $\text{RExp}_c^r U$ of L_2 would become higher than that of L_1 :

$$\begin{aligned} \text{RExp}_{c(\cdot \mid \text{Ticket } 1-11)}^r U(L_1) &= 1000 \\ \text{RExp}_{c(\cdot \mid \text{Ticket } 1-11)}^r U(L_2) &= 0 + \left(\frac{10}{11}\right)^2 \cdot (2000 - 0) \approx 1653 \end{aligned}$$

outcome gets better (i.e. as i becomes larger). In other words, the scaling factors will satisfy

$$\frac{r(p_i) - r(p_{i+1})}{p_i - p_{i+1}} > \frac{r(p_j) - r(p_{j+1})}{p_j - p_{j+1}}$$

whenever $j > i$. Since the only constraint on the p_k is that $1 > p_k > p_l > 0$ whenever $k < l$, this condition is equivalent to the claim that r is convex by a standard result [reference?? reference Buchak for other arguments?].

⁸Risk-seeking agents can be associated with *concave* risk functions, and can hence also be represented by the theory; but our focus will be on risk-avoidant agents.

For this agent, then, gathering the evidence (GATHER) is equivalent to taking L_2 , and not gathering the evidence (AVOID) is equivalent to taking L_1 . But since $\text{RExp}_c^r U(L_1) > \text{RExp}_c^r U(L_2)$, we also have that $\text{RExp}_c^r U(\text{AVOID}) > \text{RExp}_c^r U(\text{GATHER})$. Look-I therefore fails in exactly the way described earlier.

In what follows we will assume that if risk-avoidance is rationally permissible, then the kind of risk-avoidance described by Buchak's theory, with some choice of a risk-avoidant r , is rationally permissible.⁹ Part of the motivation is dialectical: Buchak is the one who suggests that advocates of risk-aversion retreat to Look-E given the failures of Look-I. Part of the motivation is principled: Buchak's theory is elegant, intuitive, and well-developed, and thus a leading contender for what rational risk-aversion might look like. And part of the motivation is pragmatic. An important advantage Buchak's theory has over other theories of risk-aversion is that it neatly separates out an agent's attitude to risk from her beliefs on the one hand, and her utilities on the other.¹⁰ This feature makes Buchak's theory particularly well-suited to studying the interaction between risk-aversion and *epistemic* rationality: because the risk-profile is separated from the utilities, we can easily study the theory's predictions when 'what matters' isn't fixed by the agent's own desires; and because the agent's attitude to risk is separated from her beliefs, we can use a familiar epistemology (namely: Bayesianism) when we do so.¹¹

3 Epistemic Rationality of Actions

According to risk-avoidant theories, instrumental rationality can sometimes require agents to avoid relevant and free evidence. However, as Buchak emphasizes, this result is consistent with claiming that epistemic rationality might still always require agents to gather such evidence. This would make rational risk-avoidance compatible with recognizing *some* sense in which gathering evidence is always a good idea.

Such a response presupposes that norms of epistemic rationality apply to actions such as the gathering or avoiding of evidence. It isn't obvious that they do: assessments of epistemic rationality are most at home when applied to doxastic states such as beliefs and credences, or to belief-producing procedures such as inference to the best explanation or conditionalization. Equally, however, it isn't obvious that they do not. When we say that it is irresponsible to consult only one kind of source, or wise to forego a quick but ambiguous test in favour of running a more thorough analysis later, we seem to engage in some form of epistemic evaluation. We will thus grant the presupposition, and assume that actions can be described as 'epistemically rational' (or 'commendable' or 'responsible' if those sound better to you).

Epistemic rationality, as applied to actions, is plausibly understood in consequentialist terms: an action is epistemically rational if it promotes the epistemic

⁹We also assume that at least one such r is differentiable; but this looks innocuous.

¹⁰Buchak (2013, pp.34-47, 53-56)

¹¹Though we will weaken some of the Bayesian assumptions regarding updating in 7.

goods of having true beliefs and accurate credences. This is not to say that epistemic rationality is consequentialist in all domains. Several philosophers have recently argued against the idea that epistemic rationality, as applied to *beliefs* or *credences*, is consequentialist. Suppose, for example, that your overall evidence weakly points towards p , but that someone powerful guarantees that you'll later find out for sure whether p if and only if you now disregard the inconclusive information and adopt a credence of 0.5. Then a naive consequentialist approach to the rationality of beliefs will make the counter-intuitive prediction that it would be epistemically rational for you to disregard your evidence and believe p to degree 0.5; for the loss in accuracy you're likely to incur by ignoring your current evidence is outweighed by the gain in accuracy when you later conform your credence to the conclusive evidence which will then be available. Perhaps more sophisticated forms of consequentialism can avoid such conclusions.¹² But it's worth noting that the analogous consequence of consequentialism about the epistemic rationality of *actions* is actually very intuitive. If I know that running a quick first-pass test into whether p will destroy the only sample and thus prevent me from carrying out a more conclusive analysis in the future, it seems epistemically quite admirable for me to refuse. While it may be weird to factor in the consequences on the accuracy of one's later beliefs, or beliefs in other propositions, when wondering whether to believe p , it is perfectly normal to do so when deciding what evidence to gather.¹³ Reasons to doubt consequentialism elsewhere in epistemology thus do not apply when the evaluations concern actions.

We will thus assume that the epistemic rationality of gathering and avoiding evidence can be understood in terms of its anticipated effects on matters of epistemic value, specifically the accuracy of one's credences.¹⁴ However, we will side-step a different controversial issue: how the accuracy of credences regarding different propositions contributes to the overall accuracy of the agent's credences.¹⁵ It is natural to think that not all propositions contribute equally (that's why we should investigate ambitious scientific theories instead of counting blades of grass); but also that there is a large amount of incommensurability

¹²Greaves (2013), Berker (2013), Caie (2013) and Carr (ms) argue that they can't. Konek and Levinstein (forthcoming) are more optimistic.

¹³Note that in cases where there are such consequences, gathering the evidence has a non-negligible epistemic cost; so we will put them aside when discussing Look-E.

¹⁴Buchak seems sympathetic to this assumption; in defending the suggestion that we can retain Look-E while rejecting Look-I, she writes that "what you have to do in connection with maximizing instrumental value is not necessarily constrained by what you have reason to do in connection with maximizing epistemic value" Buchak (2010, p.105), which suggests that epistemic rationality is a matter of pursuing epistemic value. Buchak does argue at some length that epistemic demands are not to be reduced to instrumental ones, since such a reduction is incompatible with the fact that epistemic demands are categorical. But, even if this is a good reason to reject such a reduction, it is no reason to reject a consequentialist approach to epistemic norms; otherwise the fact that moral demands are categorical would be reason to reject consequentialist theories in ethics.

¹⁵As the examples we're about to discuss bring out, this issue is particularly pressing if we want to use accuracy evaluations to determine the rationality of actions; Greaves (2013), among others, argues that it does not arise if we evaluate only the rationality of beliefs or belief-revision procedures.

(that's why it's fine to make slow progress on arcane issues in philosophy instead of reading Wikipedia all day). Fortunately, we can side-step all these issues by considering only how your action affects your accuracy on a single proposition, which the evidence in question is supposed to bear on. In this simple case, accuracy and value plausibly coincide.

Before saying more about how accuracy, and hence epistemic value, is measured, we need to discuss how exactly epistemic value determines the epistemic rationality of gathering evidence. If we assume risk-neutral decision theory, it is natural to think that an action is epistemically rational only if it maximizes *expected* epistemic value. But if we're instead considering Buchak's risk-sensitive theory, it's at least as natural to say that an agent's action is epistemically rational only if it maximizes *risk-weighted expected* epistemic value, where the risk-weightings are determined by the agent's risk function.¹⁶

There are at least two reasons we suggest that epistemic rationality should be determined by taking *risk-sensitive* expected values. The first is that this fits better with the picture painted by risk-sensitive decision theory. This theory departs from classical decision theory precisely because it maintains that a rational agent's attitude towards risk should not be understood as a generalization about her goals (e.g. that she has a concave utility function) but instead as a psychological feature that determines *how she rationally pursues* her goals.¹⁷ But if that's true, then it's natural to think that the agent's rational attitude to risk should be 'held fixed' even when the agent's subjective utilities are swapped out for epistemic values to determine the agent's epistemic obligations. The second consideration is that a major attraction of consequentialist theories quite generally is their prediction that epistemic norms require nothing 'more' than instrumental rationality once an agent's desires align with what is epistemically valuable.¹⁸ But if we determine epistemic duties in terms of expected rather than risk-weighted expected value, we lose this pleasing convergence. For these two reasons, we think that it is natural to characterize epistemic obligations in terms of risk-weighted expected epistemic value.

We will show that, if these assumptions are granted, then risk-avoidant agents will sometimes be epistemically permitted to avoid evidence. At first sight, this may look like an obvious corollary of the failure of Look-I and the view that epistemic rationality requires us to maximize risk-weighted expected epistemic value, much like the moral permissibility of avoiding evidence is an obvious corollary of the failure of Look-I and the view that morality requires us to maximize risk-weighted expected moral value. A closer look, however, reveals that this isn't so. For we have significantly fewer 'degrees of freedom' when constructing counterexamples to Look-E than when constructing counterexamples to Look-I. This is because the epistemic utilities associated with

¹⁶Another option is to say that an action is epistemically rational if it maximizes risk-weighted expected epistemic value *according to any permissible risk-profile*. We discuss how to measure accuracy if one takes this view in footnote 26, and how much of our argument this allows us to sustain in footnote 29 and footnote 32.

¹⁷See especially Buchak (2013, p.34-36).

¹⁸[Find a reference for this in the moral case.]

learning are fixed by the agent's credences after learning the evidence, which are in turn determined by the credences used in deciding whether to gather the evidence in the first place. So we cannot simply stipulate probabilities and utilities independently. Thus, (unlike in the moral case) we cannot simply use the Allais example, stipulating that the utilities in question represent epistemic, rather than instrumental, value.

This explanation also brings us to our final piece of setup. So far, we have said nothing about how epistemic value, or accuracy, is measured; but without measuring accuracy, we cannot calculate the risk-weighted expected accuracy of gathering or avoiding evidence. Showing how to measure accuracy will be our next task; it deserves its own section, since, when working with risk-weighted expectations, our accuracy measures have to behave in a somewhat non-standard way.

4 Measuring Accuracy

We are interested in measuring the epistemic utility, or accuracy, of our agent's credence in the proposition of interest, X . We can think of this as a measure of how 'close' the credence is to the truth-value, i.e. to 1 if X is true and to 0 if X is false. However, there are many ways of measuring such proximity. The obvious absolute-difference measure, on which the distance between her credence, x , and the truth value, v , is simply $|v - x|$, and the 'closeness' between them is hence $-|v - x|$, has various drawbacks; a popular alternative is the Brier Score, which measures distance as the square of the absolute difference:¹⁹

$$\text{BS}(x, v) := -(v - x)^2.$$

Following the literature, we will not defend a particular measure, but will instead adopt some general constraints, and show that our claim holds given any measure \mathcal{A} meeting these constraints. Two of the constraints we will appeal to are uncontroversial, and require no extra comment:²⁰

- \mathcal{A} is (weakly) truth directed, i.e.
 - If $x_1 < x_2 < 1$, then $\mathcal{A}(x_1, 1) \leq \mathcal{A}(x_2, 1)$,
 - If $x_1 > x_2 > 0$, then $\mathcal{A}(x_1, 0) \leq \mathcal{A}(x_2, 0)$.

¹⁹See especially Joyce (2009), Leitgeb and Pettigrew (2010a) and Pettigrew (2016, ch. 4). Most of the literature discusses measures of *inaccuracy*. But these are straightforwardly adapted, by letting the accuracy be the negative of the inaccuracy. For our purposes, this slight mathematical inelegance is worthwhile, since it preserves the connection between 'higher value' and 'better outcome', thus making it easier to apply Buchak's theory. However, nothing of substance depends on it.

²⁰See e.g. Joyce (2009). A further constraint, continuity, will be met by our specific example, but will not be required for the general theorem. Other standard constraints, such as those known as 'Normality' and 'Separability', concern the relationship between the accuracy of particular beliefs and the accuracy of the overall belief state; since we are focusing only on the accuracy of a single belief, these constraints have no bearing on our discussion.

- \mathcal{A} is 0/1 symmetric,²¹ i.e.

$$- \mathcal{A}(x, 1) = \mathcal{A}(1 - x, 0).$$

Our final condition, however, does require additional discussion – not only because it is more controversial, but also because it will need to be spelled out in a slightly non-standard way.

Many authors writing on the connection between credence and truth have argued that epistemically rational agents should be immodest: they should regard their own credence in X as giving the best shot at the truth, compared to any other credences.²² This thought comes in both a weak and a strong version. On the weak version, it merely requires that rational agents should regard their own credence in X as giving them *no worse* a shot at the truth than any other; on the strong version, it requires that agents should regard their own credence as giving them a *better* shot at the truth than any other.

One standard motivation is that someone who isn't immodest exhibits internal conflict. This is clearest in the case of weak immodesty: if you think some other credence has a better shot at the truth than yours, you seem divided in much the way as when you believe a contradiction. But it may also motivate strong immodesty: if you think that some other credence has just as good a shot at truth as yours does, it feels as though, in some sense, you're also taking this other credence, thus bearing rival attitudes towards X .

Another standard thought in motivation is that modest states are problematically unstable. If you fail to be strongly immodest, you have no epistemic reason to stick with your beliefs: if the opportunity arises, you might as well abandon it for one of the alternatives you think equally good. And if you fail to even be weakly immodest, you will not only lack reason to remain, but actually have positive reason to switch. Since such changes of mind look epistemically irrational (no new evidence is required to initiate them), this again suggests that rational agents must be immodest.

Most of the literature treats these arguments as establishing strong, rather than merely weak, immodesty;²³ and we will follow this trend in our initial discussion. However, the arguments for strong immodesty are less conclusive than those for weak immodesty. Moreover, there may be special reasons why strong immodesty is too strict a constraint when allowing for risk-avoidance; for while there are ways of measuring the accuracy of *one's attitude to a single proposition* that vindicate strong immodesty even in our risk-avoidant setting, [REFERENCE OMITTED] shows that there are no ways of measuring the accuracy of an *entire credal state* that vindicate strong immodesty in this framework.

²¹ In fact, our theorem only requires the much weaker constraint that $\mathcal{A}(x, 0) \leq \mathcal{A}(x, 1)$ when $x > 1/2$ and $\mathcal{A}(x, 0) \geq \mathcal{A}(x, 1)$ when $x < 1/2$.

²² We take talk of 'best shot at the truth' from Horowitz (forthcoming). Other sympathetic discussions include Lewis (1971), Oddie (1997), Greaves and Wallace (2006), Gibbard (2008), and Joyce (2009).

²³ See e.g. Oddie (1997), Greaves and Wallace (2006), and Joyce (2009). Maher (2002) and Gibbard (2008) raise doubts about strict immodesty, understood as a constraint on accuracy measures.

Now, for the reasons discussed in 3, the former kind of measure is more obviously relevant to our project, and so this is not a conclusive reason to reject strong immodesty; but it may, nonetheless, temper one's enthusiasm for this requirement. So, after presenting the initial argument in a way that presupposes strong immodesty, in section 6 we will consider how it fares given only weak immodesty.

We can leverage immodesty (of either kind) into a precise constraint on \mathcal{A} given two things: a sufficient condition for a distribution being rational, and an account of 'how good a shot' a distribution gives to a particular credence in X .

Plausibly, a sufficient condition for the rationality of a distribution is that it's probabilistically coherent – any coherent distribution could be rational, given the right sort of evidence.²⁴ So any such distribution should give itself a (or the) best shot at the truth.

Understanding how good a shot a distribution c gives to a credence y is slightly subtler. If we were working in the risk-neutral framework, it would be natural to identify a credence's shot at the truth with its expected accuracy. We can then guarantee that all probabilistic agents are weakly immodest by requiring \mathcal{A} to be *weakly proper*: $\text{Exp}_c \mathcal{A}(c(X)) \geq \text{Exp}_c \mathcal{A}(y)$. And we can guarantee that all probabilistic agents are strongly immodest by requiring \mathcal{A} to be *strictly proper*: for every probabilistic c and $y \neq c(X)$, $\text{Exp}_c \mathcal{A}(c(X)) > \text{Exp}_c \mathcal{A}(y)$.

However, we are assuming that some rational agents are risk sensitive and instead calculate expectations taking risk into account. And it seems more natural to say that, when an agent has credal state c and risk profile r , the credences in X she considers to give her a 'best' shot at the truth are the ones that maximize *risk-weighted* expected accuracy. We can then say that \mathcal{A} is *weakly r -proper* if for all probabilistic credence functions, c , $\text{RExp}_c^r \mathcal{A}(y)$ is maximal at $y = c(X)$; and that \mathcal{A} is *strictly r -proper* if this maximum is unique.

If we are to have a single measure of accuracy that is appropriate for every rational agent, and which also preserves weak or strong immodesty, it must be weakly or strictly r -proper for every rational risk-profile r . Unfortunately, if there are multiple rational risk profiles, this condition is very hard to meet. For if r_1 and r_2 are distinct, continuous risk profiles, then \mathcal{A} cannot be both strictly r_1 -proper and (even weakly) r_2 -proper.²⁵ This shows that we cannot have an accuracy measure such that every probabilistic agent comes out as strongly

²⁴One might worry that this is not a plausible assumption for risk avoidant agents; for example, it might never be rational for such agents to adopt an extremal credence, since doing so is very risky. However, our result goes through even if the inaccuracy measure doesn't render *all* coherent distributions immodest. All we really require is that there is some range such that the inaccuracy measure renders immodest the probability functions c defined on $\{X, \neg X\}$ for which $c(X)$ falls within the range. And the rationality of risk aversion provides no challenge to this weaker assumption. (Thanks to Jason Konek and Richard Pettigrew for discussion on this point.)

²⁵Pettigrew (2016, Section 16.4) shows that for all r , \mathcal{A} , x and y ,

$$\text{RExp}_{c_x}^r \mathcal{A}(y) = \text{Exp}_{c_{r(x)}} \mathcal{A}(y)$$

where c_x denotes the probability function on the algebra $\{X, \neg X\}$ where $c_x(X) = x$, $c_x(\neg X) = 1 - x$. So, if we pick $x_1 \neq x_2$ with $r_1(x_1) = r_2(x_2)$ (which is possible by the intermediate

immodest when epistemic utility is determined by that measure; and that even the requirement that the accuracy measure be such that every probabilistic agent comes out as weakly immodest relative to it is implausibly restrictive.

This means that if we want both immodesty and risk-avoidance we need to give up on a certain picture of epistemic utility – a picture on which there is a single quantity, epistemic value, and that every rational agent is epistemically required to pursue that quantity. The most similar (‘objectivist’) picture holds that for each risk-profile there is a unique value quantity which all agents with that risk-profile are epistemically required to maximize. (Compare: we may reject ‘single value quantity’ consequentialism in moral philosophy on the grounds that their histories and personal relationships require different people to prioritize differently the well-being of various people; but we might still say that there is a single function, which determines for each kind of history and personal circumstance a unique quantity which people with that history and circumstances are morally required to pursue.) We would then place as a constraint on this function that it associates a risk-profile r with an r -proper accuracy measure, thus ensuring immodesty.

An alternative, ‘subjectivist’, picture maintains that agents choose a particular accuracy measure, subject to certain constraints, and are then epistemically required to pursue accuracy as evaluated by that measure. We could then ensure immodesty by imposing as one of the constraints that the chosen accuracy measure be r -proper relative to the agent’s risk-profile. A third alternative is a ‘permissivist’ picture. On this account agents are not required to pursue epistemic value as measured by any particular measure. Rather, given an agent’s risk profile there are a number of legitimate epistemic utility measures; and an agent is epistemically permitted to perform any action that maximizes risk-weighted epistemic utility on some legitimate accuracy measure or other. We could then impose immodesty by maintaining that \mathcal{A} is a legitimate accuracy measure, for an agent with risk profile r , only if \mathcal{A} is r -proper.²⁶

We have now encountered all the constraints on accuracy measures that we will appeal to. We think they are all reasonable. We know that all three can be jointly met in the case of a risk-neutral agent (for whom $r(x) = x$), since in that

value theorem using the assumption that r_1 and r_2 are continuous), then

$$\begin{aligned} \text{RExp}_{c_{x_1}}^{r_1} \mathcal{A}(y) &= \text{Exp}_{c_{r_1(x_1)}} \mathcal{A}(y) \\ &= \text{Exp}_{c_{r_2(x_2)}} \mathcal{A}(y) \\ &= \text{RExp}_{c_{x_2}}^{r_2} \mathcal{A}(y) \end{aligned}$$

which is maximised at $y = x_2 \neq x_1$ if \mathcal{A} is weakly r_2 -proper. But this means that \mathcal{A} is not strictly r_1 -proper.

²⁶We mentioned in footnote 16 that epistemic rationality might not care which particular risk-profile an agent has adopted, but only about which ones are rational. This generates a view we will call *super-permissivism*: instead of aggregating verdicts across legitimate accuracy measures like the permissivist, the super-permissivist aggregates them across legitimate *standards of evaluation* consisting of both a risk-profile and an accuracy measure. We can still impose immodesty by maintaining that $\langle r, \mathcal{A} \rangle$ is a legitimate standard of evaluation only if \mathcal{A} is r -proper. As discussed in footnote 29 and footnote 32, much of our argument goes through on this super-permissivist picture as well.

case popular measures such as the Brier Score will do. However, we also know that this measure won't work for risk-avoidant agents, since no single measure is proper relative to two different risk-profiles. Fortunately, there are suitable measures even for risk-avoidant profiles; relative to $r(x) = x^2$, for example, the following measure meets all three conditions:²⁷

$$\text{AltBS}(x, v) := \frac{-(v - x)^2}{\max\{x, 1 - x\}}$$

Now that we have explained the conditions on our accuracy measures, and have provided an example of an accuracy measure suitable for a risk-avoidant agent, we can finally address our central question: whether epistemic rationality ever permits avoiding evidence.

5 Avoiding Evidence

The principle we are interested in is:

Look-E. In ideal cases, one is epistemically required to gather the evidence.

We will show that, if Buchak's risk avoidant agents are rational, and if accuracy measures should be strictly r -proper, there are counterexamples to this principle. In fact, we will show that there are even counterexamples to the weaker

Weak Look-E. In ideal cases, one is always epistemically *permitted* to gather the evidence.

Before we do this, it's worth reporting the status of these principles in classical decision theory. Suppose that c is a probability function that makes a piece of evidence, E , relevant to our proposition of interest, X . Suppose further that it's certain that, were our agent to learn whether E , she would conditionalize on what she learns. Then we get the following result regarding the expected accuracy achieved by finding out whether E (GATHER) or not finding out whether E (AVOID):

Theorem 5.1. *If \mathcal{A} is strictly proper, then $\text{Exp}_c \mathcal{A}(\text{GATHER}) > \text{Exp}_c \mathcal{A}(\text{AVOID})$.*²⁸

We thus vindicate Look-E, and hence also Weak Look-E, if rationality requires risk-neutrality, regardless of whether we opt for an 'objectivist', 'subjectivist', or 'permissivist' picture of the accuracy measures.

So much for the classical theory; on to Buchak. We will begin with a specific example, assuming the 'subjectivist' picture sketched above. We have a rational

²⁷For a graph of this function and the proof that this satisfies the desiderata, see appendix A. In unpublished work Ben Levinstein develops a method which can be applied to almost any risk-profile r to yield an inaccuracy measure \mathcal{A} that has all the properties we require (relative to r).

²⁸Oddie (1997)

agent whose risk-profile is given by $r(x) = x^2$, and whose accuracy measure is AltBS. X is the proposition of interest; $c(X) = 0.7$, $c(X|\neg E) = 0.6$, and $c(X|E) = 0.8$. Our agent is still certain that, if she finds out whether E , she will conditionalize on what she learns, so that her credence in X will move to either .8 (if E) or .6 (if $\neg E$). She is also certain that, if she declines to find out whether E , her credence in X will stay at .7. So the epistemic utility achieved by finding out whether E (GATHER) or not finding out whether E (AVOID) depends on X and E as follows:

	$E \wedge \neg X$	$\neg E \wedge \neg X$	$\neg E \wedge X$	$E \wedge X$
AVOID	AltBS(0.7, 0)	AltBS(0.7, 0)	AltBS(0.7, 1)	AltBS(0.7, 1)
GATHER	AltBS(0.8, 0)	AltBS(0.6, 0)	AltBS(0.6, 1)	AltBS(0.8, 1)

Crunching the numbers, and including information about the probability of the four different underlying states, this becomes

probability	$E \wedge \neg X$	$\neg E \wedge \neg X$	$\neg E \wedge X$	$E \wedge X$
	0.1	0.2	0.3	0.4
AVOID	-0.7	-0.7	-0.129	-0.129
GATHER	-0.8	-0.6	-0.4	-0.05

We can then calculate the risk-weighted expected accuracies of the two options, as usual considering the outcomes from worst to best:

$$\begin{aligned}
 \text{RExp}_c^r \mathcal{A}(\text{AVOID}) &= -0.7 + (0.3 + 0.4)^2 \cdot (-0.129 - -0.7) \\
 &\approx -0.420 \\
 \text{RExp}_c^r \mathcal{A}(\text{GATHER}) &= -0.8 \\
 &\quad + (0.2 + 0.3 + 0.4)^2 \cdot (-0.6 - -0.8) \\
 &\quad + (0.3 + 0.4)^2 \cdot (-0.4 - -0.6) \\
 &\quad + (0.4)^2 \cdot (-0.05 - -0.4) \\
 &= -0.484
 \end{aligned}$$

We thus get that $\text{RExp}_c^r \mathcal{A}(\text{GATHER}) < \text{RExp}_c^r \mathcal{A}(\text{AVOID})$. So even though E is relevant and the investigation is cost-free our agent is not epistemically required to find out whether E . On the contrary, our agent is epistemically required to avoid finding out whether E , since this a better means to achieving accurate opinions. Look-E and Weak Look-E both fail.

This counterexample relies on the choice of risk-profile $r(x) = x^2$ and accuracy measure AltBS. An advocate of Buchak's decision theory could in principle reject either of those: perhaps $r(x) = x^2$ is not a rational risk-profile after all, or perhaps there are constraints on accuracy other than those mentioned in section 4 which mean that AltBS is not an acceptable accuracy measure even for an agent with that risk-profile. However, a response along these lines cannot succeed. For the example just given is simply an illustration of a far more general result (proved in appendix B):

Theorem 5.2. *Suppose that r is a risk-profile that is differentiable and risk-avoidant. Then there is some prior credence distribution c over the algebra generated by X and E such that, for every \mathcal{A} which is weakly truth-directed, 0/1-symmetric, and strictly r -proper, $\text{RExp}_c^r \mathcal{A}(\text{GATHER}) < \text{RExp}_c^r \mathcal{A}(\text{AVOID})$*

This means that, regardless of which exact risk-profile our agent has, and regardless which accuracy measure she chooses within the constraints we have defended, if she is risk-avoidant then there will always be ideal cases in which she is epistemically required to avoid the evidence. Weak Look-E and Look-E thus fail.

As just formulated the point presupposes the ‘subjectivist’ picture of accuracy measures. But, clearly, the theorem similarly shows that the ‘objectivist’ picture leads to the failure of Weak Look-E and Look-E. The theorem even establishes the failure of both Look-E and Weak Look-E on the permissivist picture, which requires us to avoid the evidence only if all the legitimate measures recommend avoiding. This is because the theorem states not only that for every accuracy measure there will be an example in which it recommends avoiding the evidence (which is sufficient to refute Weak Look-E on the subjectivist and objectivist pictures and Look-E on the permissivist picture), but also that there will be a single example in which all the relevant accuracy measures recommend against gathering the evidence (which is needed to refute Weak Look-E on the permissivist picture).²⁹

All the pictures of how accuracy considerations determine epistemic obligations thus lead to the conclusion that Look-E and Weak Look-E both fail if Buchak-style risk-avoidance is rational.

Buchak offered an informal explanation of the failure of Look-I in terms of the risk of receiving *instrumentally misleading* evidence. A similar explanation is available here. Evidence can be *epistemically misleading* moving an agent’s credence further away from the truth. Suppose that E is evidence *for* X and $\neg E$ is evidence *against* X , as in our examples. Then in the state $E \wedge \neg X$, by learning E the agent will increase her credence in the false proposition X , thereby becoming less accurate. Similarly in $\neg E \wedge X$, by gathering the evidence the agent will decrease her credence in the true proposition X . In both $E \wedge \neg X$ and $\neg E \wedge X$ then, the evidence is misleading, and the epistemic utility decreases by gathering the evidence. Drawn out in our table, we have:

	$E \wedge \neg X$	$\neg E \wedge \neg X$	$\neg E \wedge X$	$E \wedge X$
AVOID	$\mathcal{A}(c(X), 0)$ \vee	$\mathcal{A}(c(X), 0)$ \wedge	$\mathcal{A}(c(X), 1)$ \vee	$\mathcal{A}(c(X), 1)$ \wedge
GATHER	$\mathcal{A}(c(X E), 0)$	$\mathcal{A}(c(X \neg E), 0)$	$\mathcal{A}(c(X E), 1)$	$\mathcal{A}(c(X \neg E), 1)$

Of course, E and $\neg E$ can also be truth-guiding, as in the two scenarios in which X and E have the same truth value. For the risk-neutral agent, this potential

²⁹Note, however, that on the super-permissivism described in footnote 26, according to which standpoints of evaluation can be relevant even though they don’t share the agent’s risk-profile, we will only get failures of Look-E and none of Weak Look-E; for risk-neutrality is surely rationally permissible, and by theorem 5.1, the standpoints of evaluation that use a risk-neutral risk-profile will always recommend gathering the evidence.

truth-guiding benefit always outweighs the risk of being misled; but for the risk-avoidant agent, we've shown that it might not.³⁰

This completes our main argument. Before closing, however, we will show that much of our conclusion can be established even if we significantly weaken two of the assumptions: that accuracy measures must be *strictly* r -proper and that rational agents update by conditionalization.

6 Weak Propriety

The previous section discussed the status of Look-E and Weak Look-E, on the assumption that accuracy measures must be strictly r -proper. However, the arguments for strict r -propriety are less compelling than those for weak r -propriety. And, as we mentioned in section 4, there are specific worries that in the risk-sensitive setting strict r -propriety becomes impossibly demanding. It's thus worth considering what happens to Look-E and Weak Look-E if we assume only that accuracy measures should be *weakly* r -proper.

Let us begin, again, with the risk-neutral case. The standard result is easily adapted to show

Theorem 6.1. *If \mathcal{A} is weakly proper $\text{Exp}_c \mathcal{A}(\text{GATHER}) \geq \text{Exp}_c \mathcal{A}(\text{AVOID})$*

This establishes Weak Look-E on any of our three pictures: gathering evidence might not always be required, but at least it is never forbidden.

Moreover, there is at least one account on which this, combined with theorem 5.1, could be used to defend Look-E. When discussing the 'permissivist' picture, we aggregated recommendations across different measures by universal agreement: an option is required only when it uniquely maximizes value according to every measure. However, we could instead have aggregated recommendations by a kind of 'pareto dominance': an option is required if it (perhaps non-uniquely) maximizes value according to every measure, and uniquely maximizes value according to at least one measure. Combined with the plausible thought that *some* legitimate measures are strictly proper, theorem 5.1 and theorem 6.1 then show that gathering the evidence is always required after all.

So much for the risk-neutral case. What about the risk-sensitive one? Here too it's easy to adapt the result to include weakly r -proper measures:

Theorem 6.2. *Suppose that r is a risk-profile that is differentiable and risk-averse. Then there is some prior credence distribution c over the algebra gen-*

³⁰One might wonder what, in general, characterizes the cases where Look-E fails. An extension of our result shows that Look-E fails whenever $c(X|E)$ and $c(X|\neg E)$ are close enough to $c(X)$, i.e. the evidence in question is relatively uninformative (and that $c(X) \neq 1/2$). (See corollary B.9, which also gives the particular form of the 'close enough'.)

This is different to the instrumental case; there Buchak (2010, Appendix C) observes that Look-I fails when the evidence is reasonable informative but not decisive. But this difference makes sense: in the instrumental case, the evidence needs to be reasonably informative to affect the agent's actions; whereas in the epistemic case, any E that is evidentially relevant at all will affect the agent's credences.

erated by X and E such that, for every \mathcal{A} which is weakly truth-directed, 0/1-symmetric,

- if \mathcal{A} is weakly r -proper, then $\text{RExp}_c^r \mathcal{A}(\text{GATHER}) \leq \text{RExp}_c^r \mathcal{A}(\text{AVOID})$
- if \mathcal{A} is strictly r -proper, then $\text{RExp}_c^r \mathcal{A}(\text{GATHER}) < \text{RExp}_c^r \mathcal{A}(\text{AVOID})$

Clearly, this is enough to show that Look-E fails: we are not always required to gather evidence. However, since, on most pictures, the retreat to weak propriety left us with no argument for Look-E even in the risk-neutral case, this does not conclusively establish any connection between risk-aversion and epistemic norms on evidence gathering.

On some pictures of the accuracy measures, we can still demonstrate such a connection. Most saliently, given the discussion of the risk-neutral case, we can do better on a version of the permissivist picture on which verdicts are aggregated via pareto dominance. For then Look-E and Weak Look-E both hold in the risk neutral case. And Look-E and Weak Look-E both fail in the risk averse case on the plausible assumption that for some risk averse risk profile there is at least one legitimate strictly r -proper accuracy measure.³¹

We can also get a connection on the subjectivist picture, on which agents can pick any accuracy measure meeting certain constraints. For even if we allow that strict r -propriety is not a constraint on how agents make their choice, plausibly, the constraints imposed still *allow* for some strictly r -proper measures. The original 5.2 will thus be enough to show that those agents who have such strictly r -proper measures will be required to avoid evidence. So we get that the rationality of risk-avoidance requires us to reject Weak Look-E, while we were able to show that this principle is true if we're required to be risk-neutral.

On the other pictures, we are restricted to less decisive plausibility arguments. Given the objectivist picture, for example, what we need to get violations of Weak Look-E is that some risk-averse risk-profile is associated with an accuracy measure for which the inequality is strict. If any such risk-profile is associated with a strictly r -proper measure, that is sufficient for the failure of Weak Look-E. And even if every risk-profile is associated only with weakly r -proper measure, it would still be surprising if in all the cases produced by 6.2, the weak inequality held because GATHER and AVOID have exactly the same expected inaccuracy. This does not conclusively establish a connection between risk-avoidance and evidence-avoidance, further constraints on accuracy measures would be required to do that. But it does make it look overwhelmingly plausible that such a connection does obtain. We can say something similar for a version of the permissivist picture which aggregates by universal verdicts rather than pareto dominance.³²

³¹This assumption is plausible even in light of [REFERENCE OMITTED]'s impossibility result mentioned above. That result shows that there can't be measures which are strictly proper when they evaluate entire credal states. However, there can clearly be measures which are (a) weakly proper in their evaluations of entire credal states and (b) strictly proper in their evaluation of credences in a particular proposition X . (Just consider the 'global' measure which cares only about your accuracy with respect to X .) And measures which are strictly

		Weak Look-E		Look-E	
		risk neutral	risk aware	risk neutral	risk aware
Subjectivist		yes	no	no	no
Objectivist		yes	probably no	no	no
Permissivist	Universal	yes	probably no	no	no
	Pareto	yes	no	yes	no

Table 1: Look Principles given only Weak Propriety

Based on these results, summarized in table 1, we conclude that, regardless of one's picture, our theorem provides at least a plausibility argument for a substantial connection between risk-avoidance and evidence-avoidance, even if strict propriety is rejected as a constraint on accuracy measures.

7 Updating by Conditionalization

To calculate the expected accuracy associated with finding out whether E , we need to know what credence the agent will adopt in response to what she learns. Since we are concerned with ideal cases, this means that we need to know what credence the agent should adopt in response to what she learns. In 5, we assumed that this is the agent's prior credence conditional on what she learns – i.e. that rational agents update by conditionalization.

Is this assumption plausible when working with risk-sensitive agents? Several of the ordinary arguments for conditionalization (e.g. that it yields intuitively plausible verdicts in many examples, or, following van Fraassen's Muddy Venn Diagram model (1989), that to update agents should simply renormalize their distribution after the possibilities inconsistent with what they learned have been eliminated) are independent of decision theory, and are thus as legitimate as ever. Others however, such as Dutch-book arguments or arguments from expected accuracy, are potentially more problematic, since they tend to assume that the preferences of rational agents are determined by expected (rather than risk-weighted expected) values.³³ We don't want to oversell the significance of

proper relative to X are all we need here.

³²What about the super-permissivism discussed in footnote 26 and footnote 29? This, too, comes in two versions, depending on whether we aggregate by dominance or by universal agreement. If we aggregate by dominance, and risk-neutrality is rationally required, theorem 5.1 and theorem 6.1 combine to establish Look-E; while if risk-avoidance is rational, theorem 6.2 shows that Look-E fails (assuming always that *some* strictly r -proper measures are acceptable). So the rationality of risk-aversion still makes a difference to Look-E (though, as in the case of strict propriety discussed in footnote 29, it makes no difference to Weak Look-E). If we aggregate by universal agreement, by contrast, we have no definite connection between risk-aversion and our principles. For Look-E will now fail even if risk-neutrality is rationally required, given the acceptability of weakly proper accuracy measures; and Weak Look-E will continue to be true, for the same reason, even if risk-aversion is rationally permissible.

³³By equating credences with betting prices, Lewis's (1999) dutch book argument assumes

this, since (1) these arguments are independently controversial (assuming, for example, that pragmatic efficacy or consequentialist reasoning is probative when it comes to the epistemic rationality of *updating rules*) and (2) it may be possible to recover these arguments for conditionalization even in a risk-sensitive setting.³⁴ But even with these reservations in mind, it's worth asking how essential the assumption that rational agents update by conditionalization is to our argument.

The answer is: not very. For theorem 5.2 holds given just a few quite general assumptions about our agent's update procedure.³⁵ Since the only cases we need are the simple ones involving just propositions X and E , we can think of an update procedure as a pair of functions $(f_E, f_{\neg E})$, which take as inputs a probability distribution c over the algebra generated by X and E , and output a real number in $[0, 1]$ which is the credence in X the agent is supposed to update to upon receiving E and $\neg E$ respectively.³⁶ Using this notation, the necessary assumptions are as follows:

1. f_E and $f_{\neg E}$ are continuous with respect to changes in the probability of the atoms.³⁷
2. If $c(E) = 1$ then $f_E(c) = c(X)$; and if $c(\neg E) = 1$ then $f_{\neg E}(c) = c(X)$.
3. $f_E(c) = c(X)$ only if $c(X) = c(X|E)$ (i.e. c makes X and E independent); and similarly for $f_{\neg E}$.
4. There is some t such that if $c(X)$, $c(X|E)$, and $c(X|\neg E)$ are all $> t$ then $f_E(c)$ and $f_{\neg E}(c)$ are both $> 1/2$.

Conditions 1. and 2. are extremely plausible: they say, respectively, that arbitrarily small changes in the input distribution shouldn't make for large changes in the output, and that learning something of which one is already certain shouldn't change one's views. Conditions 3. and 4. are also quite plausible, saying, respectively, that conditionally relevant evidence is relevant and that when X is highly probable regardless of what's the case with E , learning E shouldn't make one think X improbable. While plausible, however, the latter two thoughts both assume that conditional probabilities place substantial constraints on the update function; so someone who rejects conditionalization as fundamentally mistaken (rather than merely wrong in details) might reject one

that agents evaluate bets based on their expected payoffs. The accuracy arguments of Greaves and Wallace (2006), Leitgeb and Pettigrew (2010b), Easwaran (2013), and Pettigrew (2016, ch.4) assume that rationality requires us to update so as to maximize expected accuracy. The argument in Briggs and Pettigrew (ms) assumes only that rational agents don't use accuracy dominated belief-revision procedures, so is not subject to this problem; however, the formal result assumes that inaccuracy measures are strictly proper, and doesn't hold for the strictly r -proper measures we associate with risk-avoidant agents.

³⁴We defend this claim with respect to accuracy arguments in [REFERENCE OMITTED].

³⁵For the proof, see appendix B.

³⁶We can allow that f_E is undefined when $c(E) = 0$ and $f_{\neg E}$ is undefined when $c(\neg E) = 0$.

³⁷Perhaps a risk-aware update function can be discontinuous at particular points, perhaps because of changes in the orderings. This will not affect the result as we only need continuity as $c(X)$ converges to a fixed $c(X|E)$ (when $c(X|\neg E)$ is also fixed), or to a fixed $c(X|\neg E)$.

or both of them. So it's worth highlighting what role they play in the proof. Some condition like 3. is needed to guarantee that the update function considers a reasonably wide range of evidence to be 'potentially relevant', and hence doesn't merely validate Look-E in a vacuous way. And condition 4. is needed only to ensure that there is some range of cases in which we know how the various possible outcomes of learning the evidence are to be ordered. It is thus far from clear that denying 3. would allow one to vindicate Look-E in an interesting way – or that denying 4. would allow one to vindicate it at all.

Conditions 1.-4. are thus fair assumptions to make about the update procedure. This is not to say that they couldn't possibly be rejected.³⁸ But it does mean that the burden is not on us to show that risk avoidant agents should still update by conditionalization, but rather on our opponent to offer some reasons for thinking that risk avoidant agents should update in a radically unfamiliar – and intuitively rather odd-looking – way.

8 Conclusion

Can it be rational to avoid evidence, even when gathering it would cost you nothing, you expect it to be relevant, and you know that you would process it rationally? It's well-known that if risk avoidance is rational then it can be *instrumentally* rational to do so. But, as Buchak (2010) observes, there's more to life than instrumental rationality: there's epistemic rationality as well. What we have shown is that this does not threaten the connection; for if risk avoidance is rational, then evidence avoidance can also be epistemically rational.

More specifically, we have shown that if the kind of risk-avoidance defended by Buchak (2010, 2013) is rational, then avoiding evidence can be epistemically rational. Our main argument, presented in section 5, relies on five main assumptions: first, that people are epistemically required to pursue epistemic utility; second, that for rationally risk-avoidant agents, this amounts to maximizing risk-weighted expected epistemic utility; third, that the epistemic utility of an action is the accuracy of the anticipated resultant credences; fourth, that inaccuracy is measured according to a strictly *r*-proper accuracy measure; and, fifth, that risk-averse agents who know they will be rational expect to conditionalize on whatever evidence they encounter.

The first three assumptions are natural starting points; in section 3, we gave some brief motivations for adopting them, and distinguished them from stronger claims which might be problematic. We defended a strong version of the fourth assumption in section 4 by appeal to the thought that rational agents regard their credences as giving them the best shot at the truth; and we showed, in section 6, that much of our argument can be run from a much

³⁸Are they met by the policies motivated by risk-sensitive versions of expected accuracy arguments? In [REFERENCE OMITTED], we argue that the best ways of adapting these arguments still supports conditionalization, so that conditions 1.-4. are met. But on other ways of adapting the arguments, they support a revision procedure which is custom-made to vindicate Look-E. So one can read our theorem as showing that these revision procedures must violate at least one of our constraints.

weaker version of the fourth assumption. Finally, in section 7, we showed that the fifth assumption, too, was unnecessarily strong: the argument goes through for any update procedure satisfying a few fairly weak constraints.

We thus conclude that there is a serious tension between risk avoidance and the claim that, in ideal cases, we should always gather the evidence. But we leave open what to ultimately make of that conclusion: whether it is a reason to reject the rationality of risk avoidance, or a reason to embrace the epistemic rationality of avoiding relevant and freely available information.

References

- Maurice Allais. Fondements d'une théorie positive des choix comportant un risque et critique des postulats et axiomes de l'Ecole Americaine. *Econométrie*, pages 257–332, 1953.
- Selim Berker. Epistemic teleology and the seperateness of propositions. *Philosophical Review*, 122:337–93, 2013.
- Rachael Briggs and Richard Pettigrew. An accuracy-dominance argument for conditionalization. May 27, 2015., ms.
- Lara Buchak. Instrumental rationality, epistemic rationality, and evidence-gathering. *Philosophical Perspectives*, 24(1):85–120, 2010.
- Lara Buchak. *Risk and Rationality*. Oxford University Press Oxford, 2013.
- Michael Caie. Rational probabilistic incoherence. *Philosophical Review*, 122: 527–575, 2013.
- Jennifer Carr. Epistemic utility theory and the aim of belief. May 27, 2015., ms.
- Kenny Easwaran. Expected accuracy supports conditionalization – and conglomerability and reflection. *Philosophy of Science*, 80(1):119–142, 2013.
- Bas C. Van Fraassen. *Laws and Symmetry*. Oxford University Press, 1989.
- Allen Gibbard. Rational credence and the value of truth. Oxford University Press, 2008.
- Irving J Good. On the principle of total evidence. *The British Journal for the Philosophy of Science*, 17(4):319–321, 1967.
- Hilary Greaves. Epistemic decision theory. *Mind*, 122(488):915–952, 2013.
- Hilary Greaves and David Wallace. Justifying conditionalization: Conditionalization maximizes expected epistemic utility. *Mind*, 115(459):607–632, 2006.
- Sophie Horowitz. Accuracy and educated guesses. forthcoming.

- James Joyce. Accuracy and coherence: Prospects for an alethic epistemology of partial beliefs. Springer, 2009.
- Jason Konek and Ben Levinstein. The foundations of epistemic decision theory. *Mind*, forthcoming.
- Hannes Leitgeb and Richard Pettigrew. An objective justification of bayesianism I: Measuring inaccuracy. *Philosophy of Science*, 77(2):201–235, 2010a.
- Hannes Leitgeb and Richard Pettigrew. An objective justification of bayesianism II: The consequences of minimizing inaccuracy. *Philosophy of Science*, 77(2): 236–272, 2010b.
- David Lewis. Immodest inductive methods. *Philosophy of Science*, 38:54–63, 1971.
- David Lewis. Why conditionalize? Cambridge UP, 1999.
- Patrick Maher. Joyce’s argument for probabilism. *Philosophy of Science*, 69: 73–81, 2002.
- Graham Oddie. Conditionalization, cogency, and cognitive value. *British Journal for the Philosophy of Science*, 48:533–541, 1997.
- Richard Pettigrew. *Accuracy and the Laws of Credence*. Oxford University Press Oxford, 2016.
- A Wayne Roberts and Dale E Varberg. *Convex functions*, volume 57. Academic Press, 1974.
- Jan van Tiel. *Convex analysis*. John Wiley, 1984.
- Peter Wakker. Nonexpected utility as aversion of information. *Journal of Behavioral Decision Making*, 1:169–175, 1988.
- Stephen Willard. *General topology*. Courier Corporation, 1970.

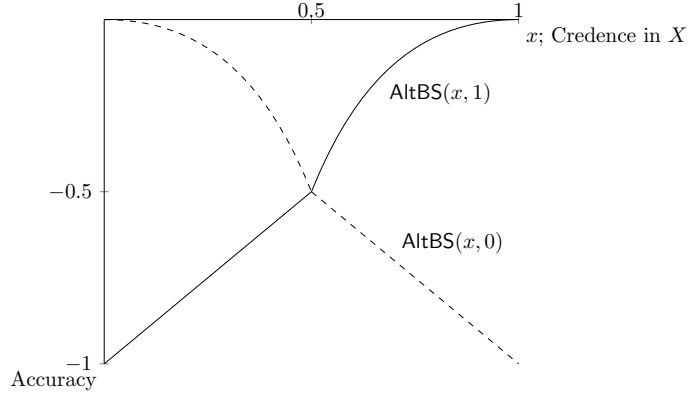
A Proof that AltBS satisfies our desiderata

Proposition A.1. *AltBS is truth-directed, 0/1 symmetric and continuous.*

Proof. Observe that

$$\text{AltBS}(y, v) := \begin{cases} y & y \geq 1/2, v = 0 \\ (1-y)^2/y & y \geq 1/2, v = 1 \\ y^2/(1-y) & y < 1/2, v = 0 \\ 1-y & y < 1/2, v = 1 \end{cases}$$

which is displayed in fig. 1 and the properties are all quite easy to check. \square

Figure 1: $\text{AltBS}(x, \cdot)$.

Proposition A.2. *AltBS is strictly r -proper.*

Proof. Assuming truth-directedness and 0/1-symmetry of \mathcal{A} ,

$$\text{RExp}_c^r \mathcal{A}(y) = \begin{cases} \mathcal{A}(y, 0) + r(c(X)) \cdot (\mathcal{A}(y, 1) - \mathcal{A}(y, 0)) & \text{if } y \geq 1/2 \\ \mathcal{A}(y, 1) + r(c(\neg X)) \cdot (\mathcal{A}(y, 0) - \mathcal{A}(y, 1)) & \text{if } y < 1/2 \end{cases}$$

The definition is piecewise because X being true is the ‘good’ case when $y \geq 1/2$, but the ‘bad’ case when $y < 1/2$.

Suppose $y \geq 1/2$. Then

$$\begin{aligned} \text{RExp}_c^r \text{AltBS}(y) &= \text{AltBS}(y, 0) + r(c(X)) \cdot (\text{AltBS}(y, 1) - \text{AltBS}(y, 0)) \\ &= y + c(X)^2 \cdot \left(\frac{(1-y)^2}{y} - y \right) \end{aligned}$$

$$\frac{\partial}{\partial y} \text{RExp}_c^r \text{AltBS}(y) = 1 - \frac{c(X)^2}{y^2}$$

so we have:

If	and	then, $\frac{\partial}{\partial y} \text{RExp}_c^r \text{AltBS}(y)$ is	so $\text{RExp}_c^r \text{AltBS}(y)$ is
$y > 1/2$	$y > c(X)$	> 0	increasing
$y > 1/2$	$y < c(X)$	< 0	decreasing

Suppose now that $y < 1/2$. Then

$$\begin{aligned} \text{RExp}_c^r \text{AltBS}(y) &= \text{AltBS}(y, 1) + c(\neg X)^2 \cdot (\text{AltBS}(y, 0) - \text{AltBS}(y, 1)) \\ &= (1-y) + c(\neg X)^2 \cdot \left(\frac{y^2}{1-y} - (1-y) \right) \end{aligned}$$

$$\frac{\partial}{\partial y} \text{RExp}_c^r \text{AltBS}(y) = 1 - \frac{c(\neg X)^2}{(1-y)^2} = 1 - \frac{(1-c(X))^2}{(1-y)^2}$$

so we have:

If	and	then, $\frac{\partial}{\partial y} \text{RExp}_c^r \text{AltBS}(y)$ is	so $\text{RExp}_c^r \text{AltBS}(y)$ is
$y < 1/2$	$y > c(X)$	> 0	increasing
$y < 1/2$	$y < c(X)$	< 0	decreasing

So we can put these two tables together into one:

If	then, $\text{RExp}_c^r \text{AltBS}(y)$ is
$y > c(X)$	increasing
$y < c(X)$	decreasing

This information tells us that there is a unique minimum at $y = c(X)$. \square

B Proof that there is always some credal state where one shouldn't Gather

We will now show that for any well-behaved and risk-avoidant r , and associated accuracy measure \mathcal{A} , there will be cases in which the expected accuracy of not gathering a relevant and cost-free piece of evidence exceeds that of gathering it.

To support our claims in section 7, we will not assume that the agent updates by conditionalization but allow her to use an arbitrary updating procedure so long as it satisfies some minimal criteria. This update procedure is given by a pair of functions $(f_E, f_{\neg E})$ which are each functions taking a prior credal states to a credal state after learning.

Our final theorem is the following:

Theorem B.1. *Suppose the following hold:*

- $r : [0, 1] \rightarrow [0, 1]$ is:
 - a risk function, i.e.
 - * increasing, and
 - * $r(0) = 0, r(1) = 1$.
 - differentiable³⁹
 - (strictly) convex, i.e. $r(\frac{x+y}{2}) < \frac{r(x)+r(y)}{2}$.
- $f_E, f_{\neg E}$ are functions from probability functions c over the algebra generated by X and E to $[0, 1]$ (though f_E may be undefined when $c(E) = 0$, and $f_{\neg E}$ when $c(\neg E) = 0$) such that:
 - f_E and $f_{\neg E}$ are continuous with respect to changes in the probability of the atoms.⁴⁰

³⁹ The argument would almost certainly go through with weaker differentiability assumptions, for example the existence of one-sided derivatives and that the right sided derivative is larger than the left-sided derivative, which follows from convexity (Tiel (1984, Theorem 1.6)). In that case, though the argument would have to be slightly modified.

⁴⁰ Though the weaker condition mentioned in footnote 37 is sufficient.

- If $c(E) = 1$ then $f_E(c) = c(X)$; and if $c(\neg E) = 1$ then $f_{\neg E}(c) = c(X)$.
- If $f_E(c) = c(X)$ then $c(X) = c(X|E) = c(X|\neg E)$ (i.e. c makes X and E independent); and similarly for $f_{\neg E}$.
- There is some $t > 1/2$ such that if both $c(X|E)$ and $c(X|\neg E)$ are $> t$ then $f_E(c)$ and $f_{\neg E}(c)$ are both $> 1/2$.

Then there is some probability distribution c over X and E such that for each \mathcal{A} where:

- $\mathcal{A} : [0, 1] \times \{0, 1\} \rightarrow \mathbb{R}$ is:
 - weakly truth directed, i.e.
 - * If $x_1 < x_2 < 1$, then $\mathcal{A}(x_1, 1) \leq \mathcal{A}(x_2, 1)$,
 - * If $x_1 > x_2 > 0$, then $\mathcal{A}(x_1, 0) \leq \mathcal{A}(x_2, 0)$.
 - 0/1-symmetric, i.e.⁴¹
 - * $\mathcal{A}(x, 0) = \mathcal{A}(1 - x, 1)$.
 - r -proper, i.e.⁴²
 - * For all $x \neq y$, $\text{RExp}_{c_x}^r \mathcal{A}(x) >/\geq \text{RExp}_{c_x}^r \mathcal{A}(y)$.
(With \geq or $>$ depending on whether it is weak or strict propriety.)

we have:

- If the r -propriety is weak, then $\text{RExp}_c^r \mathcal{A}(\text{GATHER}) \leq \text{RExp}_c^r \mathcal{A}(\text{AVOID})$
- If the r -propriety is strict, then $\text{RExp}_c^r \mathcal{A}(\text{GATHER}) < \text{RExp}_c^r \mathcal{A}(\text{AVOID})$

Proof. This follows immediately from propositions B.3 and B.5 and lemma B.2. \square

Since the update functions f_E and $f_{\neg E}$ corresponding to conditonalization have the required properties, theorems 5.2 and 6.2 are special cases of this result.

To prove this theorem we start by dealing with a special case: We note that if the update rule sometimes leads to the same value whether the agent learns E or $\neg E$, then our result is easy:

Lemma B.2. Suppose there is some credence function where $f_E(c) = f_{\neg E}(c) \neq c(X)$ then for all r -proper \mathcal{A} ,

- If the r -propriety is weak (for c), then $\text{RExp}_c^r \mathcal{A}(\text{GATHER}) \geq \text{RExp}_c^r \mathcal{A}(\text{AVOID})$

⁴¹The weaker condition discussed in footnote 21 is sufficient

⁴²The weaker condition discussed in footnote 24 is sufficient for the main result, provided the interval in question includes values greater than the t mentioned in the fourth condition on the update functions. The main result could be strengthened to say that there is a c for which $c(X)$ falls in the interval. Corollary B.9 would have to be restricted to cases where all the values fall into the interval.

- If the r -propriety is strong (for c), then $\text{RExp}_c^r \mathcal{A}(\text{GATHER}) > \text{RExp}_c^r \mathcal{A}(\text{AVOID})$

Proof. For such a c ,

$$\text{RExp}_c^r \mathcal{A}(\text{GATHER}) = \text{RExp}_c^r \mathcal{A}(f_E(c))$$

and

$$\text{RExp}_c^r \mathcal{A}(\text{AVOID}) = \text{RExp}_c^r \mathcal{A}(c(X))$$

so the result holds by r -propriety. \square

We start our main result by giving sufficient condition for the failure of Look-E:

Proposition B.3. *Suppose $\mathcal{A}, f_E, f_{\neg E}$ satisfy the conditions in theorem B.1.⁴³ Suppose $0 < c(E) < 1$ and one of the following conditions holds:⁴⁴*

1. $1/2 \leq f_{\neg E}(c) < f_E(c) \leq 1$ and

$$\frac{r(c(X)) - r(c(E \wedge X))}{r(c(\neg E) + c(E \wedge X)) - r(c(E \wedge X))} \geq r(f_E(c)).$$

2. $1/2 \leq f_E(c) < f_{\neg E}(c) \leq 1$ and

$$\frac{r(c(X)) - r(c(\neg E \wedge X))}{r(c(E) + c(\neg E \wedge X)) - r(c(\neg E \wedge X))} \geq r(f_{\neg E}(c)).$$

We then have:

- If the r -propriety is weak, then $\text{RExp}_c^r \mathcal{A}(\text{GATHER}) \leq \text{RExp}_c^r \mathcal{A}(\text{AVOID})$
- If the r -propriety is strict, then $\text{RExp}_c^r \mathcal{A}(\text{GATHER}) < \text{RExp}_c^r \mathcal{A}(\text{AVOID})$

Proof. We will only present the proof in the case where the first condition holds. The other case is completely analogous. We will use the shorthand $x = c(X)$, $y_E = f_E(c)$ and $y_{\neg E} = f_{\neg E}(c)$.

⁴³In fact no assumptions on f_E or $f_{\neg E}$ are actually required for this result. Although we have used the assumption that $f_E(c) \neq c(X)$ that isn't actually essential because one can instead use the strict propriety assumption at eq. (5).

⁴⁴One can also add the following conditions, but they won't play a role in our proof of the final theorem:

3. $0 \leq f_E(c) < f_{\neg E}(c) \leq 1/2$ and

$$\frac{r(c(\neg X)) - r(c(E \wedge \neg X))}{r(c(\neg E) + c(E \wedge \neg X)) - r(c(E \wedge \neg X))} \geq r(1 - f_E(c)).$$

4. $0 \leq f_{\neg E}(c) < f_E(c) \leq 1/2$ and

$$\frac{r(c(\neg X)) - r(c(\neg E \wedge \neg X))}{r(c(E) + c(\neg E \wedge \neg X)) - r(c(\neg E \wedge \neg X))} \geq r(1 - f_{\neg E}(c)).$$

Given this assumption of the orderings $1/2 \leq y_{\neg E} < y_E \leq 1$ and the assumptions of truth-directedness and 0/1-symmetry we see that the states are ordered:

$$\text{GATHER: } E \wedge \neg X \prec \neg E \wedge \neg X \prec \neg E \wedge X \prec E \wedge X$$

We can then calculate the risk-weighted expected utility of GATHER:

$$\begin{aligned} \text{RExp}_c^r \mathcal{A}(\text{GATHER}) &= \mathcal{A}(y_E, 0) \\ &+ r(c(\neg E \wedge \neg X) + c(\neg E \wedge X) + c(E \wedge X)) \cdot (\mathcal{A}(y_{\neg E}, 0) - \mathcal{A}(y_E, 0)) \\ &+ r(c(\neg E \wedge X) + c(E \wedge X)) \cdot (\mathcal{A}(y_{\neg E}, 1) - \mathcal{A}(y_{\neg E}, 0)) \\ &+ r(c(E \wedge X)) \cdot (\mathcal{A}(y_E, 1) - \mathcal{A}(y_{\neg E}, 1)) \end{aligned}$$

We can also observe that

$$\text{RExp}_c^r \mathcal{A}(\text{AVOID}) = \text{RExp}_c^r \mathcal{A}(c(X))$$

If we assume that \mathcal{A} is strictly r -proper, we thus have:

$$\begin{aligned} \text{RExp}_c^r \mathcal{A}(\text{AVOID}) &= \text{RExp}_c^r \mathcal{A}(c(X)) \\ &> \text{RExp}_c^r \mathcal{A}(y_E) \\ &= \mathcal{A}(y_E, 0) + r(c(X)) \cdot (\mathcal{A}(y_E, 1) - \mathcal{A}(y_E, 0)) \end{aligned}$$

And if \mathcal{A} is weakly r -proper, then we have the weak inequality:

$$\text{RExp}_c^r \mathcal{A}(\text{AVOID}) \geq \mathcal{A}(y_E, 0) + r(c(X)) \cdot (\mathcal{A}(y_E, 1) - \mathcal{A}(y_E, 0))$$

For both the strict and weak propriety results it thus suffices to show that:

$$\begin{aligned} &\mathcal{A}(y_E, 0) + r(c(X)) \cdot (\mathcal{A}(y_E, 1) - \mathcal{A}(y_E, 0)) \\ \geq &\mathcal{A}(y_E, 0) \\ &+ r(c(\neg E \wedge \neg X) + c(\neg E \wedge X) + c(E \wedge X)) \cdot (\mathcal{A}(y_{\neg E}, 0) - \mathcal{A}(y_E, 0)) \\ &+ r(c(\neg E \wedge X) + c(E \wedge X)) \cdot (\mathcal{A}(y_{\neg E}, 1) - \mathcal{A}(y_{\neg E}, 0)) \\ &+ r(c(E \wedge X)) \cdot (\mathcal{A}(y_E, 1) - \mathcal{A}(y_{\neg E}, 1)) \end{aligned} \tag{1}$$

For shorthand, let

$$a_1 = c(\neg E \wedge \neg X) + c(\neg E \wedge X) + c(E \wedge X) \tag{2}$$

$$a_2 = c(E \wedge X) \tag{3}$$

We can now rearrange eq. (1) to get that it is sufficient to show that:

$$\begin{aligned} &(r(a_1) - r(x)) \cdot \mathcal{A}(y_E, 0) + (r(x) - r(a_2)) \cdot \mathcal{A}(y_E, 1) \\ \geq &(r(a_1) - r(x)) \cdot \mathcal{A}(y_{\neg E}, 0) + (r(x) - r(a_2)) \cdot \mathcal{A}(y_{\neg E}, 1) \end{aligned}$$

Which can be further rearranged to:

$$\begin{aligned} & \mathcal{A}(y_E, 0) + \frac{r(x) - r(a_2)}{r(a_1) - r(a_2)} \cdot (\mathcal{A}(y_E, 1) - \mathcal{A}(y_E, 0)) \\ \geq & \mathcal{A}(y_{-E}, 0) + \frac{r(x) - r(a_2)}{r(a_1) - r(a_2)} \cdot (\mathcal{A}(y_{-E}, 1) - \mathcal{A}(y_{-E}, 0)) \end{aligned} \quad (4)$$

This now has exactly the form of a risk-weighted expectation, so we can now use our (at least weak) propriety assumption again. Now we use it at a credence function c_{y_E} defined on $\{X, \neg X\}$ with $c_{y_E}(X) = y_E$ and $c_{y_E}(\neg X) = 1 - y_E$:

$$\text{RExp}_{c_{y_E}}^r \mathcal{A}(y_E) \geq \text{RExp}_{c_{y_E}}^r \mathcal{A}(y_{-E}), \quad (5)$$

I.e.:

$$\begin{aligned} & \mathcal{A}(y_E, 0) + r(y_E) \cdot (\mathcal{A}(y_E, 1) - \mathcal{A}(y_E, 0)) \\ \geq & \mathcal{A}(y_{-E}, 0) + r(y_E) \cdot (\mathcal{A}(y_{-E}, 1) - \mathcal{A}(y_{-E}, 0)) \end{aligned}$$

One of the assumptions in the statement of the theorem is item 1, which we can now write in our notation as:

$$\frac{r(x) - r(a_2)}{r(a_1) - r(a_2)} \geq r(y_E)$$

Given this and the assumption about the orderings of utilities, which implies

$$\mathcal{A}(y_E, 1) - \mathcal{A}(y_E, 0) \geq \mathcal{A}(y_{-E}, 1) - \mathcal{A}(y_{-E}, 0) \geq 0$$

we get that

$$\begin{aligned} & \mathcal{A}(y_E, 0) + \frac{r(x) - r(a_2)}{r(a_1) - r(a_2)} \cdot (\mathcal{A}(y_E, 1) - \mathcal{A}(y_E, 0)) \\ \geq & \mathcal{A}(y_{-E}, 0) + \frac{r(x) - r(a_2)}{r(a_1) - r(a_2)} \cdot (\mathcal{A}(y_{-E}, 1) - \mathcal{A}(y_{-E}, 0)) \end{aligned}$$

Which is eq. (4), as required. \square

We have now got a sufficient condition for failures of Look-E. We can now show that there are such failures by showing that for differentiable risk-avoidant risk profiles, r , there will be some credences that satisfy the condition of proposition B.3.

Before proving this, we show that to exhibit credences satisfying the conditions it is sufficient to choose values for $c(X)$, $c(X|E)$ and $c(X|\neg E)$.

Lemma B.4. *If $x_{-E}, x, x_E \in [0, 1]$ and x lies strictly between x_{-E} and x_E ,⁴⁵ there is a unique probability function c defined on the algebra generated by E and X with $c(X|\neg E) = x_{-E}$, $c(X) = x$ and $c(X|E) = x_E$.*

⁴⁵I.e. either $0 \leq x_{-E} < x < x_E \leq 1$ or $0 \leq x_E < x < x_{-E} \leq 1$.

Proof. Suppose we have such $x_{\neg E}, x_E, x$. We will define the probability function on the atoms of the algebra; and then by ensuring that these sum to one (and lie in $[0, 1]$) we will know that we have in fact defined a probability function. The probabilities of the other members of the algebra are then determined by summing up these probabilities assigned to the atoms.

$$\begin{aligned} c(E \wedge X) &= x_E \cdot \frac{x - x_{\neg E}}{x_E - x_{\neg E}} \\ c(E \wedge \neg X) &= (1 - x_E) \cdot \frac{x - x_{\neg E}}{x_E - x_{\neg E}} \\ c(\neg E \wedge X) &= x_{\neg E} \cdot \frac{x_E - x}{x_E - x_{\neg E}} \\ c(\neg E \wedge \neg X) &= (1 - x_{\neg E}) \cdot \frac{x_E - x}{x_E - x_{\neg E}} \end{aligned}$$

It is easy to see that these are well defined and in $[0, 1]$. It is also a routine to verify that they sum to 1, and thus induce a probability function on the algebra generated by E and X , and that this probability function has the property that $c(X|\neg E) = x_{\neg E}$, $c(X) = x$ and $c(X|E) = x_E$. \square

We thus only need to specify values for $x_{\neg E}, x, x_E$ to exhibit a credence function satisfying the conditions of proposition B.3, and complete the proof.

Proposition B.5. *Suppose $r, \mathcal{A}, f_E, f_{\neg E}$ satisfy the constraints in theorem B.1 and:*

- *If $f_E(c) = f_{\neg E}(c)$ then they are $= c(X)$.*

Then there is some credence function c where one of the conditions in proposition B.3 holds.

Proof. We start with a useful fact about such update rules:

Lemma B.6. *One of the following holds:*

- *For all c : If $c(X) < c(X|E)$ then $f_{\neg E}(c) < f_E(c)$. Or*
- *For all c : If $c(X) < c(X|E)$ then $f_{\neg E}(c) > f_E(c)$.*

And similarly for $c(X) > c(X|E)$.

Proof. This result will use the standard theorem from topology that says that the continuous image of a connected space is connected.⁴⁶

Consider the collection of probability functions where $c(X) < c(X|E)$. We note that this is a connected set. And since f_E and $f_{\neg E}$ are continuous so is $f_E - f_{\neg E}$. So $\{f_E(c) - f_{\neg E}(c) \mid c(X) < c(X|E)\}$ is also connected by the aforementioned standard result.

Also observe that $f_E(c) - f_{\neg E}(c) \neq 0$ for all such c . This is because: Suppose $f_E(c) - f_{\neg E}(c) = 0$, i.e. $f_E(c) = f_{\neg E}(c)$. Then by the additional assumption

⁴⁶See, e.g. Willard (1970, Theorem 26.3).

in proposition B.5, $f_E(c) = c(X)$, but by the assumption in theorem B.1, then we have that $c(X) = c(X|E)$. This contradicts our assumption that we are considering credence functions where $c(X) < c(X|E)$.

So $\{f_E(c) - f_{\neg E}(c) \mid c(X) < c(X|E)\}$ is connected and $\not\equiv 0$. It must therefore be that all such credence functions have $f_E(c) - f_{\neg E}(c)$ positive, or all negative. The first would get us that $f_{\neg E}(c) < f_E(c)$ and the latter $f_{\neg E}(c) > f_E(c)$.

The same argument applies to the credence functions with $c(X) > c(X|E)$. \square

Suppose we are in the case where whenever $c(X) < c(X|E)$ then $f_{\neg E}(c) < f_E(c)$.

We thus are looking to find a c with $t < c(X|\neg E) < c(X) < c(X|E)$ (as then $1/2 \leq f_{\neg E}(c) < f_E(c)$ so we have the ordering of item 1) that satisfies:

$$\frac{r(c(X)) - r(c(E \wedge X))}{r(c(\neg E) + c(E \wedge X)) - r(c(E \wedge X))} \geq r(f_E(c)).$$

By lemma B.4 we know it suffices to choose values of $c(X|\neg E)$, $c(X|E)$ and $c(X)$.

We first pick the conditional credences, i.e. we fix some constants $d_{\neg E}$ and d_E that will stand for these conditional credences in our final credence function c . We chose them with $r(d_E) \leq d_{\neg E} < d_E$ and $t < d_{\neg E}$. This is possible because of our assumption of r being a risk avoidant risk function we have $r(x) < x$ for all $x \in (0, 1)$.

We will then consider variations of x , a variable for $c(X)$, in order to find our instance of c satisfying the conditions of proposition B.3.

Since by lemma B.4 we know that choosing any x with $d_{\neg E} < x < d_E$ determines a unique probability function we can abuse notation and consider f_E simply as a function of this variable x . We know that this is well-defined on $(d_{\neg E}, d_E)$. Moreover, we will extend $f_E(x)$ to the boundary d_E , by letting $f_E(d_E) = f_E(c)$ for c with $c(E) = 1$ and $c(X) = d_E$. The probability functions corresponding to a choice of x will be such that as $x \rightarrow d_E$, they have $c(X) - c(X|E) \rightarrow 0$, and thus the probability of E will converge to 1. So at the limit one obtains the probability function c with $c(E) = 1$ and $c(X) = d_E$. Since f_E is continuous with respect to changes in the atoms, it follows that $\lim_{x \nearrow d_E} f_E(x) = f_E(d_E)$ (and thus our extended f_E , considered as a function of just x is continuous). Moreover, by our second assumption on f_E from theorem B.1, $f_E(d_E) = d_E$.

In the proof of proposition B.3 we used a_1 and a_2 to denote the probabilities of certain events (eqs. (2) and (3)). In this theorem we will use them explicitly

as a function of x , which provides the relevant probabilities:⁴⁷

$$a_1(x) := 1 - (1 - d_E) \cdot \frac{x - d_{\neg E}}{d_E - d_{\neg E}} \quad (6)$$

$$a_2(x) := d_E \cdot \frac{x - d_{\neg E}}{d_E - d_{\neg E}} \quad (7)$$

These are well-defined for all x as the denominators are non-zero by the choice of $d_{\neg E} < d_E$.

Observe that we are therefore required to find some x where::

$$\frac{r(x) - r(a_2(x))}{r(a_1(x)) - r(a_2(x))} \geq r(f_E(x)), \quad (8)$$

as that will then provide us our credence function satisfying item 1.

Now define:⁴⁸

$$g(x) := r(x) - r(a_2(x)) - r(f_E(x)) \cdot (r(a_1(x)) - r(a_2(x)))$$

Lemma B.7. *If $g'(d_E) < 0$ then for x sufficiently close to d_E (and $x < d_E$), eq. (8) holds.*

Proof. Observe that

$$g(x) \geq 0 \iff \frac{r(x) - r(a_2(x))}{r(a_1(x)) - r(a_2(x))} \geq r(f_E(x)).$$

(so long as this right hand side is well-defined, which it is for any $x \in (d_{\neg E}, d_E)$.)

Using the fact that $a_1(d_E) = a_2(d_E) = d_E$ we can see that $g(d_E) = 0$. So since $g'(d_E) < 0$ it must be that g is decreasing at d_E and therefore for any x sufficiently close to d_E (and $< d_E$) we have

$$g(x) \geq 0. \quad \square$$

Our final lemma to prove proposition B.5 in the case where if $c(X) < c(X|E)$ then $f_{\neg E}(c) < f_E(c)$ is then:

Lemma B.8. $g'(d_E) < 0$.

⁴⁷So note that if we consider a probability function with $c(X|\neg E) = d_{\neg E}$, $c(X|E) = d_E$ and $c(X) = x$, then we have that

$$\begin{aligned} a_1(x) &= 1 - c(E \wedge \neg X) = c(\neg E \wedge \neg X) + c(\neg E \wedge X) + c(E \wedge X) \\ a_2(x) &= c(E \wedge X) \end{aligned}$$

⁴⁸Since the domain of f_E is just $(d_{\neg E}, d_E]$, the domain of g is the same. We will be considering $g'(d_E)$ so it's really one-sided derivatives that are being talked about here. But this does not affect our result because r , a_1 and a_2 are all differentiable on the whole of $(0, 1)$ and the f'_E component of g' is multiplied by zeros.

Proof. $a_1(x)$ and $a_2(x)$ (from eqs. (6) and (7)) are both differentiable with $a_2'(x) = \frac{d_E}{d_E - d_{\neg E}}$ and $a_1'(x) = -\frac{1-d_E}{d_E - d_{\neg E}}$.

So by using the chain and product rules we get:⁴⁹

$$\begin{aligned} g'(x) &= r'(x) - r'(a_2(x))a_2'(x) - r(f_E(x)) \cdot (r'(a_1(x))a_1'(x) - r'(a_2(x))a_2'(x)) \\ &\quad - r'(f_E(x)) \cdot f_E'(x) \cdot (r(a_1(x)) - r(a_2(x))) \\ &= r'(x) - r'(a_2(x))\frac{d_E}{d_E - d_{\neg E}} \\ &\quad - r(f_E(x)) \cdot \left(r'(a_1(x)) \left(-\frac{1-d_E}{d_E - d_{\neg E}} \right) - r'(a_2(x))\frac{d_E}{d_E - d_{\neg E}} \right) \\ &\quad - r'(f_E(x)) \cdot f_E'(x) \cdot (r(a_1(x)) - r(a_2(x))) \end{aligned}$$

Since $a_1(d_E) = a_2(d_E) = d_E$, the derivative at d_E is:

$$\begin{aligned} g'(d_E) &= r'(d_E) - r'(d_E)\frac{d_E}{d_E - d_{\neg E}} \\ &\quad - r(f_E(d_E)) \cdot \left(r'(d_E) \left(-\frac{1-d_E}{d_E - d_{\neg E}} \right) - r'(1-d_E)\frac{d_E}{d_E - d_{\neg E}} \right) \\ &\quad - r'(f_E(d_E)) \cdot f_E'(d_E) \cdot (r(d_E) - r(d_E)) \\ &= r'(d_E) \left(1 - \frac{d_E - r(f_E(d_E))}{d_E - d_{\neg E}} \right) \end{aligned} \tag{9}$$

Since r is convex and increasing we have $r'(d_E) > 0$.⁵⁰ We thus have

$$g'(d_E) < 0 \iff r(f_E(d_E)) < d_{\neg E}.$$

To obtain our result we thus have to show the right hand side of this. Our assumption on f_E (together with the definition of $f_E(d_E)$) ensures that $f_E(d_E) = d_E$. We chose d_E and d_E with $r(d_E) < d_{\neg E}$. So the right hand side of this holds, as required. \square

This has therefore shown that if we are in the case where if $c(X|\neg E) < c(X) < c(X|E)$ then $f_{\neg E}(c) < f_E(c)$, it must be that there is some c where condition 1 of proposition B.3 holds and thus that $\text{RExp}_c^r \mathcal{A}(\text{GATHER}) < / \leq \text{RExp}_c^r \mathcal{A}(\text{AVOID})$.

Supposing instead we were in the case where if $c(X|\neg E) < c(X) < c(X|E)$ then $f_{\neg E}(c) > f_E(c)$, we would need to redo some of this proof. However, the argument goes through pretty much similarly. The proof would then work by:

- Choose *any* $t < d_{\neg E} < d_E$.

⁴⁹ Although this refers to $f_E'(x)$ we will be taking the derivative at d_E and for that we do not need to assume this exists, roughly because it'll then be multiplied by zero. We can say this because: if $h_2(x_0) = 0$ and h_2 is differentiable, h_1 is continuous, then

$$(h_1 \cdot h_2)'(x_0) = h_1(x_0) \cdot h_2'(x_0).$$

⁵⁰ Since r is convex its derivative is increasing (Roberts and Varberg, 1974, Theorem 42B), and since r is an increasing function $r'(x) \geq 0$ for all $x \in (0, 1)$. Thus $r'(d_E) > r'(d_E - \epsilon) \geq 0$.

- Redefine $a_1(x) = 1 - (1 - d_{\neg E}) \cdot \frac{x - d_E}{d_{\neg E} - d_E}$ and $a_2(x) = x_{\neg E} \cdot \frac{x - d_{\neg E}}{d_{\neg E} - d_E}$, and the resulting g will be the same as here except using these new a_1 and a_2 , and will involve $f_{\neg E}(x)$ instead of $f_E(x)$.
- One then observes that if $g'(d_{\neg E}) > 0$ then we can get our item 2 by choosing x close enough to $d_{\neg E}$, and $d_{\neg E} < x < d_E$.
- We show $g'(d_{\neg E}) > 0$ by an analogous argument to lemma B.8. One needs to notice in the eq. (9) the denominator will be $d_{\neg E} - d_E$ which is negative and so we get that $g'(d_{\neg E}) > 0 \iff r(f_{\neg E}(d_{\neg E})) < d_E$.
- We show this inequality by observing $f_{\neg E}(d_{\neg E}) = d_{\neg E} < d_E$. \square

Suppose for simplicity the update rule is conditionalization. Then we can also get a sufficient condition:

*For all $y \in [0, 1] \setminus \{1/2\}$
For all z closer to $1/2$ than y and close enough _{y} to y*

Corollary B.9. *For all x between y and z and close enough _{y, z} to y
If c has y and z as $c(X|E), c(X|\neg E)$ and x as $c(X)$
then $\text{RExp}_c^r \mathcal{A}(\text{GATHER}) > / \geq \text{RExp}_c^r \mathcal{A}(\text{AVOID})$*

Where what counts as close enough depends on the indicated variables.

Proof. For the case where $y = c(X|E) > 1/2$ the above proof works for this just by checking what choices were made. For the other cases one can observe that much of the proof goes through *mutatus mutandis*. \square

If the update rule is not conditionalization a theorem similar to the above will still be available but it will be more complex with additional assumptions to ensure, for example, that the orderings of the states are appropriate.