

Rational Inattention with Multiple Attributes

David Walker-Jones*

University of Surrey

November 2, 2021

Abstract

This paper uses an axiomatic foundation to create a new measure for the cost of learning that allows the different attributes of the options faced by an agent to differ in their associated learning costs. The new measure maintains the tractability of Shannon’s classic measure but produces richer choice predictions and identifies a new form of informational bias significant for welfare and counterfactual analysis that is conducted with the multinomial logit model. Necessary and sufficient conditions are provided for optimal behavior under the new measure for the cost of learning.

1 Introduction

In many choice environments it is costly for agents to learn about the options that they face because it takes time and effort to acquire and process information. Understanding how agents learn in such environments is crucial for quality economic analysis because the cost of information may result in agents not acquiring all of the relevant information before making a decision. Partially informed agents do not always pick the best available option, which makes welfare analysis more challenging. Further, if what information an agent acquires changes with parameters such as price then counterfactual analysis is also made more difficult.

The standard technique for quantifying the cost of learning in models of rational inattention (RI) is Shannon Entropy ([Shannon, 1948](#); [Sims, 2003](#); [Mackowiak, Matejka, & Wiederholt, 2021](#)). Shannon Entropy has an axiomatic foundation, is grounded in the optimal coding of information,

*Special thanks to Rahul Deb for all of the support. I would also like to thank Andrew Caplin, Mark Dean, Yoram Halevy, Marcin Peski, Carolyn Pitchik, and Colin Stewart, for their helpful advice. An earlier version of this paper was previously circulated under the title “Rational Inattention and Perceptual Distance.”

and provides a tractable and flexible framework with which to study agent behavior (Shannon, 1948; Matějka & McKay, 2015; Caplin, Dean, & Leahy, 2018).

While Shannon Entropy has proven to be a valuable tool, it does have limitations in economic environments as they are not what it is designed for. It is, for instance, natural to think that some attributes of the choice environment might be more difficult to learn about than others. Shannon Entropy, however, does not allow for attributes of the choice environment to differ in their associated learning costs because it is a one parameter model for the cost of information, and thus there can only be one level of difficulty when learning. Without a mechanism to allow for what is referred to in the literature as “perceptual distance,”¹ the choice behavior predicted by Shannon Entropy can differ from observed behavior, as is discussed in Example 1 in Section 2.1, which can limit the effectiveness of Shannon Entropy in empirical settings (Dean & Neligh, 2018).

This paper proposes four axioms that are similar to Shannon’s original axioms (1948) in that they focus on the cost of answering simple questions that can be represented by partitions of the state space. Taken together, the four axioms in this paper are weaker than Shannon’s axioms (1948) because they relax Shannon’s assumption that the set of simple questions that is used, and the order in which they are answered, cannot change the expected cost of learning the state of the world. By allowing for the set of questions that is used to learn the state of the world, and the order in which they are answered, to change the expected learning cost I produce a new measure for the cost of information that I call Multi-Attribute Shannon Entropy (MASE), which allows for attributes of the choice environment to differ in their associated learning costs.

Though the axioms in this paper discuss an agent learning through simple partitions of the state space, one need not constrain the agent to learn in such a fashion, and MASE can be used to measure the cost of information in a more general model where the agent can choose to learn with any signal structure they desire, as is typical in the literature on RI. This is because MASE, like Shannon Entropy, can be viewed as a measure of total uncertainty. As such, the cost of an arbitrary signal can simply be measured as the difference between the total uncertainty before and after the signal is realized as is frequently done with Shannon Entropy in models of RI.

MASE maintains much of the desired tractability of Shannon’s classic measure when incorporated into a model of RI because this paper provides the MASE analogues of the famous necessary conditions provided by Matějka and McKay (2015) and necessary and sufficient condi-

¹If two outcomes are more difficult to differentiate between it is said that they have less perceptual distance between them.

tions provided by [Caplin et al. \(2018\)](#) for optimal behavior in RI models that use Shannon Entropy. Further, MASE predicts behavioral patterns that have been identified as problematic for Shannon Entropy. Thus, while Shannon Entropy is a one parameter model of the costs of learning, MASE provides a natural multi-parameter generalization of Shannon Entropy.

MASE is flexible enough to allow for different attributes of an economic agent’s different options to differ in their associated learning costs. MASE is thus flexible enough to, for instance, be the foundation of a model of obfuscation in which different firms choose how difficult it is for consumers to learn about the different attributes of their products, while Shannon Entropy is not even flexible enough to allow for different options to differ in their associated learning costs.

MASE also identifies an informational bias in the multinomial logit random utility (RU) model that should be considered a natural consequence of different learning costs in the same choice environment, as is demonstrated by [Example 2](#) in [Section 2.2](#). While other papers study measures of information that feature different levels of learning costs (e.g., [Hébert & Woodford, 2017](#); [Pomatto, Strack, & Tamuz, 2019](#)), this paper is the first to identify an informational bias in a standard RU model that is generated by the presence of different learning costs in the same choice environment. Unlike the informational bias identified with Shannon Entropy ([Matějka & McKay, 2015](#)), this type of informational bias cannot be identified with the average choice probabilities of the agent,² and thus presents a new challenge for welfare and counterfactual analysis.

1.1 Organization of Paper

The remainder of the paper is organized as follows: [Section 2](#) introduces Shannon Entropy, discusses models of RI, and provides motivating examples. [Section 3](#) proposes four new axioms, and uses them to develop a more flexible cost of acquiring information, MASE, which allows for attributes of the choice environment to differ in their associated learning costs. [Section 4](#) uses MASE as a benchmark with which to price inattentive information strategies in a model of RI, and discusses the resultant agent behavior, showing that much of the coveted tractability of Shannon Entropy is maintained by this paper’s generalization by establishing necessary and sufficient conditions for optimal behavior. [Section 5](#) discusses the relationship between RU models and the agent behavior found in [Section 4](#), and revisits the motivating examples from [Section 2.1](#) and [Section 2.2](#). [Section](#)

²The average choice probability of an option is the weighted average of the probabilities of it being selected in the different potential states of the world. Later in the paper this is referred to as the unconditional probability of the option being selected, as is standard in the literature, since the probability does not condition on the state of the world.

State:	ω_1	ω_2	ω_3	ω_4
Balls in State:	60 Blue & 40 Red	51 Blue & 49 Red	49 Blue & 51 Red	40 Blue & 60 Red
Probability of State:	1/4	1/4	1/4	1/4
Value of option 1:	y	y	$-y$	$-y$
Value of option 2:	0	0	0	0

[6](#) provides a literature review, and [Section 7](#) concludes.

2 Rational Inattention with Shannon Entropy

To find an introduction to models of rational inattention, and in particular models of RI that use Shannon Entropy to measure the cost of information, please see the first part of [Appendix 3](#).

Problems can occur when Shannon Entropy is applied in settings with an option that has multiple attributes that differ in their associated learning costs, as is discussed in [Example 1](#), or in settings with options that have different associated learning costs, as is discussed in [Example 2](#).

2.1 Example 1: Multiple Attributes and Problems with Predictions

[Caplin, Dean, and Leahy \(2017, p. 19\)](#) show that Shannon Entropy results in choice behavior that satisfies “invariance under compression.” That is, when Shannon Entropy is used to measure information, if there are two states of the world, ω_1 and ω_2 , across which payoffs are identical for each option ($\mathbf{v}_n(\omega_1) = \mathbf{v}_n(\omega_2) \forall n \in \mathcal{N}$), then the chance of each option being selected is the same in ω_1 and ω_2 . The invariance under compression that is predicted by Shannon Entropy is, unfortunately, not found in many settings, as is shown by the work of [Dean and Neligh \(2018\)](#). This subsection describes an environment akin to the experiments in [Dean and Neligh \(2018\)](#).

Consider the environment described in [Table 1](#) where an agent is faced with a screen that shows 100 balls, each of which is either red or blue. The agent is offered a prize that they may either accept (option 1), or reject to get a payoff of zero (option 2). The agent is told that if the majority of the balls on the screen are blue then the prize is $y \in \mathbb{R}_{++}$, and if the majority of the balls on the screen are red then the prize is $-y$. Suppose further that the agent is also told that there is a 1/4 chance of four different states of the world in which there are either 40, 49, 51, or 60 red balls.

The Shannon RI model, which imposes invariance under compression, predicts that the agent has the same chance of selecting option 1 when there are 40 red balls as when there are 49 red

balls, and that the agent has the same chance of selecting option 1 when there are 60 red balls as when there are 51 red balls. This predicted behavior is not intuitive because it should be easier for the agent to differentiate between the states that are more different (40 versus 60 red balls) than the states that are more similar (49 versus 51 red balls). One should instead expect that the chance that option 1 is selected is decreasing in the number of red balls, as is demonstrated by the experiments of [Dean and Neligh \(2018\)](#), because it should be easier to determine which color of ball constitutes the majority the more of that color ball there are.

Why does Shannon Entropy impose this type of behavior? In short, Shannon Entropy results in invariance under compression because of Shannon’s third axiom ([Shannon, 1948](#)). In the context of [Example 1](#), let $\mathcal{P} = \{\{\omega_1\}, \{\omega_2\}, \{\omega_3\}, \{\omega_4\}\}$, and $\tilde{\mathcal{P}} = \{\{\omega_1 \cup \omega_2\}, \{\omega_3 \cup \omega_4\}\}$, be two partitions of the state space. Shannon’s third axiom requires that total uncertainty about the state of the world, be equal to the uncertainty about which event in $\tilde{\mathcal{P}}$ has occurred, plus the expected amount of uncertainty that remains about which event in \mathcal{P} has occurred after which event in $\tilde{\mathcal{P}}$ occurred has been learned. This equality means that the reduction in uncertainty caused by a signal, which is the cost of the signal, is equal to the reduction in uncertainty about which event in $\tilde{\mathcal{P}}$ has occurred, plus the expected reduction in uncertainty about which event in \mathcal{P} has occurred given which event in $\tilde{\mathcal{P}}$ has occurred.

The agent is only concerned with which event in $\tilde{\mathcal{P}}$ has occurred, as this fully determines payoffs. If agent behavior is different in ω_1 compared to ω_2 , or ω_3 compared to ω_4 , so that their behavior does not satisfy invariance under compression, then the agent is, to an extent, differentiating between these states, and paying for information that does not benefit them, and their information strategy is thus not optimal.

While other information cost functions do not require that choice behavior satisfies invariance under compression ([Caplin et al., 2017](#); [Morris & Yang, 2016](#)), they lack the tractability and flexibility of Shannon Entropy,³ which limits the potential for their application. This has led to the following open question: “what workable alternative models allow for the complex behavioral patterns identified in practice?” ([Caplin et al., 2017](#), p. 2), a question that this paper attempts to answer. MASE solves the problem with predictions outlined in this example by allowing option 1 to have multiple attributes that differ in their learning cost, as is explained in [Section 5.2](#).

³Shannon Entropy has a number of mathematical properties that make it easy to use for predicting behavior in a wide range of environments.

State:	ω_1	ω_2	ω_3	ω_4
Probability of State:	1/4	1/4	1/4	1/4
Value of option 1:	H	H	L	L
Value of option 2:	H	L	H	L

2.2 Example 2: Options that Differ in Learning Costs and Biases in Fitting

If attributes vary in their learning costs then RU models are susceptible to a form of informational bias that has not previously been identified, as demonstrated by the following example. This is significant for those who wish to conduct welfare or counterfactual analysis because there are many economically significant examples where, for instance, one option is easier to learn about, as in [Example 2](#).

Consider a choice environment where an agent has two options: option 1 and option 2, which can each be of high value H , or low value $L < H$, as is described in [Table 2](#). Assume, contrary to what is possible with Shannon Entropy, that learning the value of option 1 is less costly than learning the value of option 2. For example, perhaps the agent is interested in investing in one of two businesses that are *a priori* identical except for the fact that one is local and easier to learn about, while the other is foreign and harder to learn about. It is not difficult to come up with more examples along these lines.

Because payoffs are symmetric, any knowledge about the value of option 1 has the same value to the agent as the same knowledge about option 2. Further, the cost of said information about option 1 is lower. As such, while the marginal benefit of information about option 1 or option 2 is the same, the marginal cost of information about option 1 is lower. One should thus expect research of a rational agent to be more attentive to option 1, and more cognisant of its value.

If both option 1 and option 2 have realized their high value H , one should expect that the agent is more likely to select option 1 since our intuition is that the agent should be more cognisant of option 1's high value. Similarly, if option 1 and option 2 have both realized their low value L , then one should expect that the agent is more likely to select option 2.⁴

Because of this, if an econometrician tried to deduce the two values of option 1, H_1 and L_1 , and the two values of option 2, H_2 and L_2 , using a multinomial logit regression, they would decide that H_1 is more than the true value H , and that L_1 is less than the true value L (as is shown rigorously in [Section 5](#)). Fitting thus falls prey to an informational bias, undermining the value of

⁴Our intuition is that the agent should be more cognisant of option 1's low value.

any counterfactual or welfare analysis.

This type of bias has not previously been identified in the literature on RI. [Matějka and McKay \(2015\)](#) show that fitting of multinomial logit results in the value of an option n to be biased by the (weighted) average probability of it being selected over states ω . The bias found by [Matějka and McKay \(2015\)](#) can thus be identified by examining the average probabilities of the agent selecting each option because the driving mechanism is that the cost of learning causes the agent to be biased towards options that they have a higher chance of selecting *a priori*. The bias previously found by [Matějka and McKay \(2015\)](#) is fundamentally different than the bias demonstrated in this example because their bias does not allow for an option to be over valued in some states and under valued in others, which is in contrast with our setting where option 1 is over valued when it is of high value, and is undervalued when it is of low value.

An econometrician who observes equal average choice probabilities in this setting, as is predicted by MASE, might be tempted conclude, based on the previous literature, that their analysis is not susceptible to informational biases since each option has the same chance of being selected *a priori*, and thus any counterfactual or welfare analysis that they conduct is valid. This conclusion may not be correct given the results in this paper.

RU models and RI models with Shannon Entropy can both be rejected for RI with MASE in this environment if it is possible to alter the correlation between the values of the two options while holding the marginal distributions over values fixed for each option, as is discussed in [Section 5.3](#) when this example is revisited.⁵

3 Multi-Attribute Shannon Entropy (MASE)

In this section I use axioms to develop this paper’s measure of uncertainty, which is the expected cost to the agent of perfectly observing the state of the world. The measure of uncertainty that I develop can then be used to study a rationally inattentive agent because the cost of a noisy information strategy can be taken to be the expected reduction in uncertainty, as is frequently done with Shannon Entropy in models of RI. Thus, while this paper is interested in studying an inattentive agent that only partially learns about the state of the world, this section discusses an attentive agent that perfectly observes the state of the world. Before I introduce the axioms in [Section 3.2](#) I pause to introduce some notation and terminology in [Section 3.1](#).

⁵This assertion is not difficult to show with [Theorem 2](#) and [Corollary 1](#).

3.1 Formal Setting

I am interested in an agent who is researching a measurable space (Ω, \mathcal{F}) , where Ω is a finite set of possible **states of the world** (the state space), and \mathcal{F} is the set of **events** generated by Ω (the power set of Ω). I call $\mu : \mathcal{F} \rightarrow [0, 1]$, which assigns probabilities to events, the **prior** belief.

One natural way to think about an agent learning is through a series of questions that have answers that are uniquely determined by the state of the world.⁶ How do I model such questions? A **partition** \mathcal{P} of a state space Ω is a set of more than one disjoint events in \mathcal{F} whose union is Ω .⁷

A question with multiple potential answers is thus equivalent to a partition whenever the answer to the question is deterministically determined by the state of the world. This equivalence occurs since I can simply group states of the world based on the answer to the question they produce. The words ‘question’ and ‘partition’ can thus be used interchangeably.

The simplest kind of question in this setting is a yes or no question. A yes or no question is equivalent to a **binary partition** \mathcal{P}^b of Ω , which I define as a set of two events, $\mathcal{P}^b = \{A_1, A_2\}$, such that $A_1 \cup A_2 = \Omega$, and $A_1 \cap A_2 = \emptyset$. The two phrases ‘binary partition’ and ‘yes or no question’ can thus be used interchangeably.

If $\omega \in \Omega$ is the state of the world, let the **realized event** of the partition $\mathcal{P} = \{A_1, \dots, A_m\}$ be denoted by $\mathcal{P}(\omega)$, that is $\mathcal{P}(\omega) = A_i \in \{A_1, \dots, A_m\}$ iff $\omega \in A_i$. Given a prior μ , if the agent **learns the realized events** $\mathcal{P}_1(\omega), \dots, \mathcal{P}_n(\omega)$ of a collection of partitions, their updated belief is denoted $\mu(\cdot | \cap_{i=1}^n \mathcal{P}_i(\omega))$, and is defined on states $\tilde{\omega}$ as follows:

$$\mu(\tilde{\omega} | \cap_{i=1}^n \mathcal{P}_i(\omega)) = 0 \text{ if } \tilde{\omega} \notin \cap_{i=1}^n \mathcal{P}_i(\omega), \text{ and otherwise } \mu(\tilde{\omega} | \cap_{i=1}^n \mathcal{P}_i(\omega)) = \frac{\mu(\tilde{\omega})}{\mu(\cap_{i=1}^n \mathcal{P}_i(\omega))}.$$

Given a prior μ , and some partition \mathcal{P} , let $C(\mathcal{P}, \mu) \in \mathbb{R}_+$ denote the cost of learning the realized event $\mathcal{P}(\omega)$ of \mathcal{P} . $C(\mathcal{P}, \mu)$, the cost of answering ‘What is the realized event of \mathcal{P} ?’ given the agent’s prior belief is the basic building block of this paper.

A **learning strategy**, $S = (\mathcal{P}_1, \dots, \mathcal{P}_n)$, is a list of partitions whose realized events are successively observed by the agent such that if $\mathcal{P}_i, \mathcal{P}_j \in S$, and $i \neq j$, then $\mathcal{P}_i \neq \mathcal{P}_j$. A ‘learning strategy’ is thus ‘a series of questions’ and the two phrases can be used interchangeably. If a learning strategy consists of only binary partitions, I call it a **binary learning strategy**, and denote it $S^b = (\mathcal{P}_1^b, \dots, \mathcal{P}_n^b)$. The order of the questions in a learning strategy is important, and

⁶A question’s answer is uniquely determined by the state if you know the answer if you know the state.

⁷Notice that the definition of a partition excludes trivial partitions that only contain a single event.

changing the order results in a different learning strategy. If, for instance, some questions are more costly for the agent to answer, and help to identify states that are seldom observed, then it may seem efficient for a learning strategy to leave these questions towards the end.⁸

I define $C(S, \mu)$, which is the expected cost of a learning strategy $S = (\mathcal{P}_1, \dots, \mathcal{P}_n)$ given a probability measure μ , to be the sum of the expected costs of each of the questions in S :

$$C(S, \mu) = C(\mathcal{P}_1, \mu) + \mathbb{E} \left[C(\mathcal{P}_2, \mu(\cdot | \mathcal{P}_1(\omega))) + \dots + C(\mathcal{P}_n, \mu(\cdot | \bigcap_{i=1}^{n-1} \mathcal{P}_i(\omega))) \right].$$

The definition of $C(S, \mu)$ thus imposes a form of constant marginal cost onto learning strategies because over the course of their learning strategy the agent does not fatigue, nor do they gain experience with research and become better at learning: all that matters for determining the cost of each question are the beliefs of the agent immediately before the question is answered, and not how much has previously been learned.

If B is any collection of partitions, let $\sigma(B)$ denote the σ -**algebra generated by** B , which is the smallest σ -algebra containing all the events in each of the partitions in B . Since a learning strategy S is a collection of partitions, I use $\sigma(S)$ to denote the σ -algebra generated by S .

Sometimes a single question can be as informative as several questions. I say a learning strategy S is **equivalent** to a partition \mathcal{P} if $\sigma(S) = \sigma(\mathcal{P})$.⁹ What $\sigma(S) = \sigma(\mathcal{P})$ means intuitively is that, for any prior probability measure $\mu : \mathcal{F} \rightarrow \mathbb{R}_+$, observing the answers to the series of questions in S always leads to the same posterior as observing the answer to the question ‘what is the realized event of the partition \mathcal{P} ?’, and thus, for all priors, S and \mathcal{P} provide the same amount of information.

3.2 Axioms

What form should a cost function for information take? This difficult question does not have an obvious answer, so this paper takes an axiomatic approach. The axioms in this section, like Shannon’s original axioms, focus on an attentive agent that is perfectly learning the state of the world. I want the axioms to be normatively appealing, and I find axioms about perfectly observing the state of the world to be a more intuitive, and hence easier to evaluate normatively, than axioms that focus directly on inattentive behavior and describe costs of different kinds of stochastic exper-

⁸The order of the events in a partition, in contrast, is not important, and switching the order in which the events in a partition are listed does not result in a different partition.

⁹Thus, I say that a series of questions is equivalent to a particular question if the learning strategy that represents the series of questions is equivalent to the partition that represents the particular question.

iments. The axioms make explicit the structure that is imposed on the cost function. Each axiom can be separately evaluated in different contexts, either empirically, or through introspection, to determine how appropriate it is. Further, the axioms help demonstrate to those that are familiar with Shannon’s original axioms (1948) the differences between MASE and standard Shannon Entropy.

Axiom 1 (Measurement): Given a binary partition $\mathcal{P}^b = \{A_1, A_2\}$, $C(\mathcal{P}^b, \mu)$ is determined by $\mu(A_1)$ and $\mu(A_2)$, and I can thus write $C(\mathcal{P}^b, \mu) = C(\mathcal{P}^b, \mu(A_1), \mu(A_2))$.

In plain language, [Axiom 1](#) says that the expected cost of the yes or no question represented by \mathcal{P}^b should be fully determined by the chance that the answer is yes and the chance that the answer is no. If I know the yes or no question being asked, and the chance of each of its answers, then I know the expected cost of answering the question, I do not require any additional information.

I am now going to introduce learning strategy invariance, a concept that is the central pillar of Shannon’s (1948) axioms and helps to make it explicit what I am assuming with this paper’s axioms. In general, a particular question \mathcal{P} , and an equivalent series of questions S , may produce different expected costs depending on what questions are selected, and how they are ordered in S . A given question \mathcal{P} , however, may have the peculiar property that, given any prior, all series of questions that are equivalent to it have the same expected cost, in which case I say it is learning strategy invariant. Formally, I say a partition \mathcal{P} is **learning strategy invariant**, if for each probability measure μ , the expected cost $C(S, \mu)$ is the same for every learning strategy S that is equivalent to \mathcal{P} .

In many environments there are questions that are not learning strategy invariant. Consider the environment described in [Example 2](#) in [Section 2.2](#). In this context, let $A_1 = \{\omega_1, \omega_2\}$, $A_2 = \{\omega_1, \omega_3\}$, $\mathcal{P}_1^b = \{A_1, A_1^c\}$, and $\mathcal{P}_2^b = \{A_2, A_2^c\}$. Notice that observing the realized event of \mathcal{P}_1^b is equivalent to learning the value of option 1, and observing the realized event of \mathcal{P}_2^b is equivalent to learning the value of option 2. Now, let $\mathcal{P}_3 = \{\{\omega_1\}, \{\omega_2\}, \{\omega_3\}, \{\omega_4\}\}$ denote the partition of the state space. Notice that the learning strategy $S^b = (\mathcal{P}_1^b, \mathcal{P}_2^b)$ is equivalent to \mathcal{P}_3 , because if I answer ‘What is the value of option 1?’, and then answer ‘What is the value of option 2?’, I have observed the state of the world.

Based on the discussion in [Section 2.2](#), however, I should expect that \mathcal{P}_3 may not be learning strategy invariant. Consider $\tilde{S}^b = (\mathcal{P}_2^b, \mathcal{P}_1^b)$, which is also equivalent to \mathcal{P}_3 . If the value of option 1 and option 2 were perfectly correlated, then observing the value of one of them would tell you the

value of the other. The cost of S^b would then be the cost of observing the value of option 1, which I assumed to be less than the cost of observing the value of option 2, which is then the cost of \tilde{S}^b .

A set of partitions that are certainly learning strategy invariant, in contrast, is the set of binary partitions. If \mathcal{P}^b is a binary partition, then \mathcal{P}^b is learning strategy invariant because the only learning strategy S such that $\sigma(S) = \sigma(\mathcal{P}^b)$, is $S = (\mathcal{P}^b)$. Thus, for any μ , all learning strategies S such that $\sigma(S) = \sigma(\mathcal{P}^b)$ have the same expected cost $C(S, \mu) = C(\mathcal{P}^b, \mu)$.

To begin to use the axioms I show that if \mathcal{P} is a leaning strategy invariant partition comprised of three or more events, then $C(\mathcal{P}, \mu)$ is constant with respect to permutations of the probability measure μ on \mathcal{P} . If $\mathcal{P} = \{A_1, \dots, A_m\}$ is a leaning strategy invariant partition, I say that $\tilde{\mu}$ is a **permutation** of μ on \mathcal{P} if there is a bijection $\pi : \{1, \dots, m\} \rightarrow \{1, \dots, m\}$ such that $\forall i \in \{1, \dots, m\}, \mu(A_i) = \tilde{\mu}(A_{\pi(i)})$.

Lemma 1. If a partition $\mathcal{P} = \{A_1, \dots, A_m\}$ is learning strategy invariant, with $m \geq 3$, and C satisfies [Axiom 1](#), then $C(\mathcal{P}, \mu)$ is fully determined by $\mu(A_1), \mu(A_2), \dots$, and $\mu(A_m)$, and if $\tilde{\mu}$ is a permutation of μ on \mathcal{P} , then $C(\mathcal{P}, \mu) = C(\mathcal{P}, \tilde{\mu})$.

Proofs for results in [Section 3](#) and can be found in [Appendix 1](#)

I next show that if a partition $\mathcal{P} = \{A_1, \dots, A_m\}$ is learning strategy invariant with $m \geq 3$, structure is imposed onto $C(\mathcal{P}^b, \mu)$ for all \mathcal{P}^b coarser than \mathcal{P} . I say a partition \mathcal{P} of a state space Ω is **coarser** than a partition $\tilde{\mathcal{P}}$ of the same state space Ω , if each event in \mathcal{P} corresponds to a union of events in $\tilde{\mathcal{P}}$. As it turns out, this structure is quite helpful.

Lemma 2. If a partition $\mathcal{P} = \{A_1, \dots, A_m\}$ is learning strategy invariant with $m \geq 3$, and \mathcal{P}^b is a binary partition coarser than \mathcal{P} , then if C satisfies [Axiom 1](#), then for all (p_1, p_2, p_3) such that $p_1, p_2, p_3 \in [0, 1)$ and $p_1 + p_2 + p_3 = 1$:

$$\begin{aligned} & C(\mathcal{P}^b, p_1, 1 - p_1) + (1 - p_1)C\left(\mathcal{P}^b, \frac{p_2}{p_2 + p_3}, \frac{p_3}{p_2 + p_3}\right) \\ &= C(\mathcal{P}^b, p_2, 1 - p_2) + (1 - p_2)C\left(\mathcal{P}^b, \frac{p_1}{p_1 + p_3}, \frac{p_3}{p_1 + p_3}\right) \\ &= C(\mathcal{P}^b, p_3, 1 - p_3) + (1 - p_3)C\left(\mathcal{P}^b, \frac{p_1}{p_1 + p_2}, \frac{p_2}{p_1 + p_2}\right). \end{aligned}$$

The next axiom, [Axiom 2](#), is concerned with the optimal order for questions to be asked in. In Shannon's (1948) original work he assumes that permuting the order of questions does not

change the expected cost of learning the state of the world. This might not make sense, however, if different attributes of the choice environment have different learning costs associated with them. If some attributes are less expensive to learn about then it might make sense to learn about these attributes first, as is argued in this subsection in the context of Example 2. Axiom 2 weakens Shannon's assertion and only imposes that if two binary partitions are similar enough then the order in which their realized events are learned about does not change the expected learning cost.

The most succinct and objective way to discuss a partition being similar to another one is with a product space. Consider replicating the state space three times so that the new state space is $\tilde{\Omega} = \Omega_1 \times \Omega_2 \times \Omega_3$ with $\Omega_1 = \Omega_2 = \Omega_3 = \Omega$. Suppose \mathcal{P}_1^b is a binary partition of Ω_1 , that \mathcal{P}_2^b is the equivalent binary partition of Ω_2 , and that \mathcal{P}_3^b is the equivalent binary partition of Ω_3 . Thus, \mathcal{P}_1^b , \mathcal{P}_2^b , and \mathcal{P}_3^b , are as similar as partitions can be by construction. Now, suppose that the agent knows that the answer to one of the questions, \mathcal{P}_1^b , \mathcal{P}_2^b , or \mathcal{P}_3^b , is 'yes,' while the other two have the answer 'no.'¹⁰ Denote the probability of \mathcal{P}_i^b having the answer 'yes' by $p_i \in [0, 1)$ for $i \in \{1, 2, 3\}$.¹¹ Suppose the agent begins by leaning about the realized event of \mathcal{P}_i^b . If the agent learns the answer to \mathcal{P}_i^b is 'yes' they know the state of the world and they are done learning, and if they instead learn the answer to \mathcal{P}_i^b is 'no' then their belief is updated so the probability of the answer to \mathcal{P}_j^b being 'yes' for $j \in \{1, 2, 3\} \setminus \{i\}$ is $\frac{p_j}{p_j + p_k}$, where $k \in \{1, 2, 3\} \setminus \{i, j\}$, and after they learn about the answer to \mathcal{P}_j^b they are done learning. What [Axiom 2](#) imposes is that the order in which the agent answers these ostensibly identical questions is irrelevant to their expected learning cost.

Axiom 2 (Self-Similarity): Given a binary partition \mathcal{P}^b , and a vector of probabilities (p_1, p_2, p_3) such that $p_1, p_2, p_3 \in [0, 1)$ and $p_1 + p_2 + p_3 = 1$, C is such that:

$$\begin{aligned} & C(\mathcal{P}^b, p_1, 1 - p_1) + (1 - p_1)C\left(\mathcal{P}^b, \frac{p_2}{p_2 + p_3}, \frac{p_3}{p_2 + p_3}\right) \\ &= C(\mathcal{P}^b, p_2, 1 - p_2) + (1 - p_2)C\left(\mathcal{P}^b, \frac{p_1}{p_1 + p_3}, \frac{p_3}{p_1 + p_3}\right) \\ &= C(\mathcal{P}^b, p_3, 1 - p_3) + (1 - p_3)C\left(\mathcal{P}^b, \frac{p_1}{p_1 + p_2}, \frac{p_2}{p_1 + p_2}\right). \end{aligned}$$

Next I make a very weak assumption about the continuity of the cost function on binary

¹⁰If the answer to one question does not contain information about the answers to the other questions, then assuming that the order in which they are answered does not impact expected costs is a vacuous assumption. The assumption made here is perhaps the simplest way to ensure the answer to one question provides information about the answer to the other questions.

¹¹The open upper bound on the p_i ensures the agent does not already know the state of the world.

partitions. As such, the axioms do not explicitly rule out discontinuities in the cost function, but, later results show that the cost function is continuous on binary partitions. This is because the property described in [Axiom 2](#) is only compatible with a cost function that is either continuous or discontinuous at every point for each binary partition.

Axiom 3 (Weak continuity): Given a binary partition \mathcal{P}^b , there is a probability $p \in [0, 1]$ such that C is continuous at $(p, 1 - p)$ when applied to \mathcal{P}^b .

As was alluded to, a cost function on binary partitions only satisfies [Axiom 1](#) and [Axiom 2](#) if it is either continuous everywhere or discontinuous everywhere. Thus, if a cost function on binary partitions satisfies the first three axioms, it is continuous everywhere, as is formalized by [Lemma 3](#), which further shows that the cost function is permutation invariant on binary partitions.

Lemma 3. If C satisfies [Axiom 1](#), [Axiom 2](#), and [Axiom 3](#), then for each binary partition \mathcal{P}^b , $C(\mathcal{P}^b, p, 1 - p)$ is continuous in p , and $C(\mathcal{P}^b, p, 1 - p) = C(\mathcal{P}^b, 1 - p, p)$, for each $p \in [0, 1]$.

Continuity and symmetry are not the only helpful properties imposed onto the cost function by the axioms. On binary partitions, the cost function is also non-decreasing if the chance of whichever event is less likely increases.

Lemma 4. If C satisfies [Axiom 1](#), [Axiom 2](#), and [Axiom 3](#), then for each binary partition \mathcal{P}^b , and for each $p \in [0, \frac{1}{2})$, $C(\mathcal{P}^b, p, 1 - p)$ is non-decreasing for small increases in p .

I now show that the cost of learning with a learning strategy invariant partition is dictated by Shannon Entropy.

Lemma 5. If a partition \mathcal{P} is learning strategy invariant, and C satisfies [Axiom 1](#), [Axiom 2](#), and [Axiom 3](#), then there exists a multiplier $\lambda(\mathcal{P}) \in \mathbb{R}_+$, such that for all probability measures μ : $C(\mathcal{P}, \mu) = \lambda(\mathcal{P})\mathcal{H}(\mathcal{P}, \mu)$, where \mathcal{H} is Shannon's standard measure of entropy ([1948](#)), defined in equation ([14](#)).

Underlying each learning strategy invariant partition is some attribute of the choice environment. [Shannon \(1948\)](#) imposes learning strategy invariance onto all partitions of Ω , which implies that all partitions have the same costs associated with them (there is a $\lambda > 0$ such that $\lambda(\mathcal{P}) = \lambda$ for all partitions \mathcal{P} of Ω), and so it is without loss to think of the agent as learning about a single attribute that allows them to differentiate between the different states of the world. With MASE,

in contrast, different learning strategy invariant partitions are allowed to have different costs associated with them ($\lambda(\mathcal{P})$ may differ depending on the learning strategy invariant partition \mathcal{P}), and thus it is natural to think of the agent as learning about different attributes of the choice environment depending on which attribute allows them to acquire the information at the lowest costs, as is formalized by [Theorem 1](#). This interpretation is how MASE gets its name.

Shannon’s (1948) key axiom, his third axiom, assumes that all partitions of the state space are learning strategy invariant, and further, that the cost function derived is defined for vectors of arbitrary length, even though Shannon also uses a finite state space to derive his cost function. In addition to this axiom, Shannon has two other axioms, one of which imposes continuity onto his cost function (his axiom 1), and another that deals with the cost of differentiating between a greater number of equally likely states (his axiom 2). As it turns out, there is a great deal of redundancy in Shannon’s axioms, as is demonstrated by this paper’s axioms.

As a result, Shannon’s third axiom is the only axiom that it is substantive to relax. Shannon’s second axiom does not have any impact as long as leaning with binary partitions is assumed to be costly when there is uncertainty about their realized event. Removing his first axiom only has an impact if I allow for a cost function that is discontinuous at every point when applied to a binary partition, which would render it too complex and intractable for practical application. As a result, if one wishes to generalize Shannon Entropy to achieve a more flexible but still tractable tool with which to study an environment where learning is costly, it must be Shannon’s third axiom that is weakened.

The first three axioms do not rule out that learning with a binary partition can be costless. I, however, wish to study a costly learning environment.¹²

Axiom 4 (Costly Learning): Given a binary partition \mathcal{P}^b , $C(\mathcal{P}^b, \frac{1}{2}, \frac{1}{2}) > 0$.

To ease exposition slightly, [Axiom 4](#) imposes that answering yes or no questions is costly to the agent.

3.3 Total Uncertainty

[Lemma 5](#) and [Axiom 4](#) together tells us that for each binary partition \mathcal{P}^b , there is an **associated multiplier**, $\lambda(\mathcal{P}^b) \in \mathbb{R}_{++}$, such that for all probability measures μ : $C(\mathcal{P}^b, \mu) =$

¹²Allowing for costless learning is not difficult theoretically, but it does make exposition slightly more cumbersome. It can be shown that if free information is available then it is optimal for the agent to acquire that information, and then given its realization, choose an optimal learning strategy as described by the results in this paper.

$\lambda(\mathcal{P}^b)\mathcal{H}(\mathcal{P}^b, \mu)$. Since there are a finite number of binary partitions of Ω , I can order the binary partitions by their associated multipliers. Let λ_1 denote the multiplier associated with all binary partitions, denoted $\{\mathcal{P}_i^{b,\lambda_1}\}_{i=1}^{n_1}$, with the lowest multiplier.

If the agent can always learn the state of the world by asking questions with multiplier λ_1 , then $\sigma(\{\mathcal{P}_i^{b,\lambda_1}\}_{i=1}^{n_1}) = \mathcal{F}$, and $M=1$.¹³ If not, let λ_2 denote the multiplier associated with all binary partitions, denoted $\{\mathcal{P}_i^{b,\lambda_2}\}_{i=1}^{n_2}$, with the second lowest multiplier.

If the agent can always learn the state of the world by asking questions with multipliers λ_1 or λ_2 , then $\sigma(\{\mathcal{P}_i^{b,\lambda_1}\}_{i=1}^{n_1}, \{\mathcal{P}_i^{b,\lambda_2}\}_{i=1}^{n_2}) = \mathcal{F}$, and $M = 2$. If not, let λ_3 denote the multiplier associated with all binary partitions, denoted $\{\mathcal{P}_i^{b,\lambda_3}\}_{i=1}^{n_3}$, with the third lowest multiplier.

Continue in this fashion until λ_M denote the multiplier associated with all binary partitions, denoted $\{\mathcal{P}_i^{b,\lambda_M}\}_{i=1}^{n_M}$, with the lowest multiplier such that the state of the world is always revealed when all questions with equal or lower associated multipliers are asked, that is, the lowest M such that: $\sigma(\{\mathcal{P}_i^{b,\lambda_1}\}_{i=1}^{n_1}, \dots, \{\mathcal{P}_i^{b,\lambda_M}\}_{i=1}^{n_M}) = \mathcal{F}$.

To help make the notation more compact, one can use a group of partitions to **generate** a finer partition: if $(\mathcal{P}_1, \dots, \mathcal{P}_n)$ is a group of partitions, let $\times\{\mathcal{P}_i\}_{i=1}^n$ denote the partition such that $\sigma(\times\{\mathcal{P}_i\}_{i=1}^n) = \sigma(\mathcal{P}_1, \dots, \mathcal{P}_n)$. Then, for $j \in \{1, \dots, M\}$,¹⁴ let $\mathcal{P}_{\lambda_j} = \times\{\mathcal{P}_i^{b,\lambda_j}\}_{i=1}^{n_j}$.

The partitions described in the preceding paragraphs are the foundation for the different attributes of the choice environment. More specifically, the **attributes** $\mathcal{A}_j \equiv \mathcal{P}_{\lambda_j}$ for $j \in \{1, \dots, M\}$ are just specific partitions of the state space since the different outcomes for each attribute divide the state space into events. That is, $\forall \omega \in \Omega$ the **realization of the attribute** \mathcal{A}_j is defined $\mathcal{A}_j(\omega) \equiv \mathcal{P}_{\lambda_j}(\omega) \in \mathcal{F}$.

Finally, since Ω is a partition of itself, one can, as a minor abuse of notation, let $S^b(\Omega) = \{S^b | \sigma(S^b) = \mathcal{F}\}$ denote the set of binary learning strategies such that $\sigma(S^b) = \sigma(\Omega) = \mathcal{F}$. I define the minimal cost of learning the state of the world in terms of the minimal cost of learning the state of the world with yes or no questions since this is the focus of the axioms.¹⁵

Theorem 1. If C satisfies all four axioms, then there exists attributes $\mathcal{A}_1, \dots, \mathcal{A}_M$, as defined

¹³If $M=1$, then MASE collapses to standard Shannon Entropy.

¹⁴ M is defined in the preceding paragraphs.

¹⁵There are a number of reasons to focus on learning with yes or no questions. Eye tracking analysis shows that when agents are faced with multiple options, they successively compare pairs of the options along a single attribute dimension (Noguchi & Stewart, 2014, 2018). This suggests that, in practice, agents are breaking their learning into a number of smaller queries. Further, in the psychology literature these pairwise comparisons are frequently modelled as ordinal in nature (Noguchi & Stewart, 2018), equivalent to questions with binary outcomes, e.g. ‘Is option a better than option b in dimension x ?’, instead of more complicated questions, e.g. ‘How much better is option a than option b in dimension x ?’, because findings in the field of psychophysics suggest that agents are good at discriminating stimuli, but are not good at determining the magnitude of the same stimuli (Stewart, Chater, & Brown, 2006).

above, and unique constants $0 < \lambda_1 < \dots < \lambda_M$ such that for any probability measure μ on \mathcal{F} :

$$\min_{S \in \mathcal{S}^b(\Omega)} C(S, \mu) = \lambda_1 \mathcal{H}(\mathcal{A}_1, \mu) + \mathbb{E} \left[\lambda_2 \mathcal{H}(\mathcal{A}_2, \mu(\cdot | \mathcal{A}_1(\omega))) + \dots + \lambda_M \mathcal{H}(\mathcal{A}_M, \mu(\cdot | \bigcap_{i=1}^{M-1} \mathcal{A}_i(\omega))) \right],$$

where \mathcal{H} is defined as in equation (14).

In plain language, [Theorem 1](#) says that if the cost of learning satisfies all four axioms, then the cheapest way (in expectation) to learn the state of the world always involves first learning the realization of the attribute \mathcal{A}_1 , the cheapest attribute to learn about, then learning the realization of attribute \mathcal{A}_2 , the second cheapest attribute to learn about, and continuing in this fashion until the state of the world has been realized.

[Theorem 1](#) generates the more flexible measure of uncertainty that I desired for studying inattentive behavior. If the agent starts with a prior μ , and does optimal learning that reaches a posterior $\tilde{\mu}$, then I let the cost of this inattentive research be measured as the reduction in uncertainty, as is discussed in the next section.

The \mathcal{P}_{λ_i} 's that are used to generate the attributes in [Theorem 1](#) are not unique, with the exception of \mathcal{P}_{λ_1} , and thus the attributes are not unique. The versions described in the paragraphs preceding [Theorem 1](#) are the unique coarsest partitions that could be used to define the attributes in the statement of the theorem. For $i \in \{2, \dots, M\}$, \mathcal{P}_{λ_i} could, for instance, be replaced by $\tilde{\mathcal{P}}_{\lambda_i} = \times \{\mathcal{P}_{\lambda_j}\}_{j=1}^i$ for generating \mathcal{A}_i in the statement of [Theorem 1](#), which would constitute the unique finest representation of the partitions that could be used to define the attributes.

The axiomatic derivation of the cost benchmark in this paper requires a discrete state space for the state of the world, as is the case with Shannon Entropy. If a continuous state space is desired for the state of the world, however, a measure of uncertainty for a continuous state space can be defined in an analogous manner to the measure of uncertainty defined in [Theorem 1](#) for a discrete state space, which is similar to what is done by [Shannon \(1948\)](#) to apply Shannon Entropy in a continuous setting.

4 Inattentive Learning with MASE

The following section introduces and solves a model of RI that uses MASE to measure the cost of acquiring information and establishes that MASE can be incorporated tractably into a model of RI, which is not an obvious result. Apart from the use of MASE instead of Shannon Entropy for

the measurement of uncertainty, this section follows the work of [Matějka and McKay \(2015\)](#) closely so as to aid comparison between the two models.

Given the result in [Theorem 1](#), I take the expected cost of a particular information strategy to be defined as:

$$\mathbf{C}(F(s, \omega), \mu) \equiv \mathbb{E} \left[\min_{S \in \mathcal{S}^b(\Omega)} C(S, \mu) - \min_{S \in \mathcal{S}^b(\Omega)} C(S, \mu(\cdot|s)) \right]. \quad (1)$$

A noisy information strategy reduces the total amount of uncertainty, and I thus measure the cost of such a noisy information strategy as the expected reduction in total uncertainty. This interpretation can also be applied to RI models that use Shannon Entropy to measure the cost of noisy information structures.¹⁶

The cost functions that can be defined as above with MASE are in the class of uniformly posterior-separable cost functions described by [Caplin et al. \(2017\)](#). The behavior generated in static settings by such posterior-separable cost functions has been shown to be equivalent to the behavior generated by sequential information sampling in some dynamic contexts ([Hébert & Woodford, 2017](#); [Morris & Strack, 2019](#)). In particular, [Hébert and Woodford \(2017\)](#) show that a class of static cost functions, which they call ‘neighborhood-based’ cost functions, can be micro-founded in this way. The cost functions explored in this paper that measure the reduction in MASE are a strict subset of the neighborhood-based cost functions described in their paper. Thus, the cost functions in this paper are micro-founded in two ways, directly through the axioms in this paper, and indirectly through the dynamic analysis conducted by [Hébert and Woodford \(2017\)](#). While symmetry imposes a unique set of partitions in [Example 1](#) when MASE is used, there are numerous representations that can be used when a neighborhood-based cost function is assumed. [Hébert and Woodford \(2017\)](#) suggest two ways of modelling the neighborhoods in such a setting, one of which is fitted by [Dean and Neligh \(2018\)](#), and neither of which is equivalent to the partitions suggested by MASE.

[Huettner, Boyacı, and Akçay \(2019\)](#), in turn, create an ad hoc group of cost functions that are also a generalization of Shannon Entropy, but are a strict subset of the cost functions studied in this paper that measure reduction in MASE. The cost functions studied by [Huettner et al. \(2019\)](#) allow different options to have different learning costs associated with them, but are not capable of

¹⁶Shannon Entropy is a measure of total uncertainty derived from axioms about the cost of successively learning the realized events of partitions, and in such models the cost of a noisy signal is simply taken to be the reduction in total uncertainty, as measured by Shannon Entropy.

predicting the behavior I argued was intuitive in [Example 1](#), since in [Example 1](#) their cost functions collapses to standard Shannon Entropy as they do not allow for a single option to have different attributes that vary in their associated learning costs.

4.1 Rationally Inattentive Agent’s Problem

As is introduced in [Section 3.1](#), the agent faces uncertainty described by the probability space $(\Omega, \mathcal{F}, \mu)$. Suppose that an agent that has stopped learning must make a selection from a set of **options**, denoted $\mathcal{N} = \{1, \dots, N\}$. Each option, $n \in \mathcal{N}$, in each state of the world, $\omega \in \Omega$, has a (finite) **value** to the agent $\mathbf{v}_n(\omega)$. To ease exposition, for the rest of the paper I assume $\mu(\omega) > 0 \forall \omega \in \Omega$.

The agent’s problem is to maximize the expected value of their selected option less the cost of learning. They do this by choosing an **information strategy** $F(s, \omega) \in \Delta(\mathbb{R} \times \Omega)$, which is a joint distribution between s , the observed **signal**, and the states of the world.¹⁷ The only restriction on the information strategy is that the marginal, $F(\omega) : \mathcal{F} \rightarrow \mathbb{R}_+$, must equal the prior μ . It is a property of MASE, as is true with Shannon Entropy, that if $F(s, \omega)$ is optimal, then the agent is done learning after a single signal s . After the signal is realized, the agent simply picks the action with the highest expected value:

$$a(s|F) = \arg \max_{n \in \mathcal{N}} \mathbb{E}_{F(\omega|s)}[\mathbf{v}_n(\omega)].$$

Ignoring the cost of learning momentarily, the value to the agent of receiving a signal s , which induces posterior $F(\omega|s)$, is then:

$$V(s|F) = \max_{n \in \mathcal{N}} \mathbb{E}_{F(\omega|s)}[\mathbf{v}_n(\omega)].$$

The agent’s problem is to maximize the expected value of the option they select less the cost of learning by choosing an optimal information strategy, and subsequently selecting an option based on the signal produced by their information strategy. The agent’s problem can thus be written:

$$\max_{F \in \Delta(\mathbb{R} \times \Omega)} \sum_{\omega \in \Omega} \int_s V(s|F) F(ds|\omega) \mu(\omega) - \mathbf{C}(F(s, \omega), \mu), \quad (2)$$

¹⁷The decision to allow s to be any real number is rather arbitrary. This is a richer signal space than is required in practice. It is shown that an optimal strategy only results in one of at most N different signals being observed.

$$\text{such that } \forall \omega \in \Omega : \int_s F(ds, \omega) = \mu(\omega). \quad (3)$$

The above problem is complicated and not particularly tractable, so I follow [Matějka and McKay \(2015\)](#) and re-write this problem directly in terms of the choice probabilities of the agent. This process requires the development of some new notation. Define $\mathcal{S}(n|F) = \{s \in \mathbb{R} : F(s) > 0, a(s|F) = n\}$, to be the set of signals that result in the agent selecting option n . Next, define the chance of option n being selected conditional on the state of the world to be:

$$\Pr(n|\omega) = \int_{s \in \mathcal{S}(n|F)} F(ds|\omega), \quad (4)$$

and for event $A \in \mathcal{F}$, define the chance of n being selected conditional on A being realized to be:

$$\Pr(n|A) = \sum_{\omega \in A} \Pr(n|\omega)\mu(\omega|A). \quad (5)$$

Define the **unconditional choice probability** of option n to be:

$$\Pr(n) = \sum_{\omega \in \Omega} \Pr(n|\omega)\mu(\omega). \quad (6)$$

Denote the collection $\{\Pr(n|\omega)\}_{n=1}^N$ by \mathbb{P} . Using this notation, I can re-write the agent's problem:

Lemma 6. Choice probabilities \mathbb{P} are the outcome of a solution to the agent's problem in (2) subject to (3) iff they solve:

$$\max_{\mathbb{P}} \sum_{n \in \mathcal{N}} \sum_{\omega \in \Omega} \mathbf{v}_n(\omega) \Pr(n|\omega) \mu(\omega) - \mathbf{C}(\mathbb{P}, \mu), \quad (7)$$

$$\text{such that: } \forall n \in \mathcal{N}, \Pr(n|\omega) \geq 0, \forall \omega \in \Omega, \quad (8)$$

$$\text{and } \sum_{n \in \mathcal{N}} \Pr(n|\omega) = 1 \forall \omega \in \Omega, \quad (9)$$

where $\mathbf{C}(\mathbb{P}, \mu)$ is as defined in [Lemma 14](#).

Proofs for results in [Section 4](#) and [Section 5](#) can be found in [Appendix 2](#)

This new problem, where the agent selects their conditional choice behavior \mathbb{P} , is substantially easier to solve than the problem where the agent picks their information strategy $F(s, \omega)$.

4.2 Behavior of a Rationally Inattentive Agent

Using [Lemma 6](#), I can establish a necessary condition for the optimal behavior of the agent with [Theorem 2](#), and then use said necessary condition to simplify the maximization problem undertaken by the agent with [Corollary 1](#).

Theorem 2.

If \mathbb{P} is the solution to (7) subject to (8) and (9), then $\forall n \in \mathcal{N}$, and $\forall \omega \in \Omega$, the probability that option n is selected in state ω satisfies:

$$\Pr(n|\omega) = \frac{\Pr(n)^{\frac{\lambda_1}{\lambda_M}} \Pr(n|\mathcal{A}_1(\omega))^{\frac{\lambda_2 - \lambda_1}{\lambda_M}} \dots \Pr(n|\cap_{i=1}^{M-1} \mathcal{A}_i(\omega))^{\frac{\lambda_M - \lambda_{M-1}}{\lambda_M}} e^{\frac{\mathbf{v}_n(\omega)}{\lambda_M}}}{\sum_{\nu \in \mathcal{N}} \Pr(\nu)^{\frac{\lambda_1}{\lambda_M}} \Pr(\nu|\mathcal{A}_1(\omega))^{\frac{\lambda_2 - \lambda_1}{\lambda_M}} \dots \Pr(\nu|\cap_{i=1}^{M-1} \mathcal{A}_i(\omega))^{\frac{\lambda_M - \lambda_{M-1}}{\lambda_M}} e^{\frac{\mathbf{v}_\nu(\omega)}{\lambda_M}}}. \quad (10)$$

Those familiar with the work of [Matějka and McKay \(2015\)](#) will recognize the above formula as the MASE analogue of [Matějka and McKay \(2015\)](#)'s Theorem 1. When there is only one attribute to learn about and $\lambda_1 = \lambda_2 = \dots = \lambda_M$, the above formula collapses to [Matějka and McKay \(2015\)](#)'s Theorem 1.

With standard Shannon Entropy, the chance that the agent selects an option thus depends only on the unconditional chances of the options being selected, and the realized values of the options. With MASE, in contrast, as the above formula indicates, the chance that the agent selects an option n in a particular state of the world ω depends on the unconditional chances of the options being selected, $\Pr(n)$, the realized values of the options $\mathbf{v}_n(\omega)$, as well as the realizations of attributes that are easier to learn about. It makes sense that when easier to observe pieces of information indicate that an option n is likely of above average value, that the agent should select option n with a higher probability, even if the above average value has not been realized. For a more complete [discussion](#) of the intuitive properties of the choice behavior described in [Theorem 2](#), please see the second part of [Appendix 3](#).

Behavior that is consistent with [Theorem 2](#) is not necessarily optimal because in many settings it is not optimal for the agent to consider all of the available options (choose them with positive probability), and though such a corner solution may be optimal, there are many corners that are

consistent with [Theorem 2](#) but are not optimal. For instance, for any $n \in \mathcal{N}$, if the agent selects n with probability one in all states of the world, then their behavior is consistent with [Theorem 2](#), but it is easy to come up with examples where this would not be optimal for any n , as is demonstrated when [Example 1](#) and [Example 2](#) are revisited in [Section 5](#).

Given conditional choice probabilities \mathbb{P} , I define the **consideration set** to be $\mathcal{C}(\mathbb{P}) \equiv \{n \in \mathcal{N} | \Pr(n) > 0\}$. I say an option n is **considered** if $\Pr(n) > 0$. This definition of a consideration set has the advantage that it can be observed in the data and fits with the definition given by [Caplin et al. \(2018\)](#).

Suppose conditional choice probabilities $\tilde{\mathbb{P}}$ that are consistent with [Theorem 2](#) and are a candidate for optimal behavior. To determine if it is in fact optimal for an option $n \in \mathcal{N}$ that is not considered under $\tilde{\mathbb{P}}$ to be considered, n needs to be compared to a representative value of the options that are being considered under $\tilde{\mathbb{P}}$ and given a score in each state of the world. The agent would do better with n in the consideration set if it scores well enough across all states of the world, in which case $\tilde{\mathbb{P}}$ is not optimal even though it is consistent with [Theorem 2](#).

In state ω , option n is compared to a weighted average of the options currently being considered, but the weight assigned to an option ν that is currently being considered depends on the unconditional probability of it being selected, as well as each of the conditional probabilities of it being selected given the realizations of the attributes $\Pr(\nu | \cap_{i=1}^j \mathcal{A}_i(\omega))$ for $j \in \{1, \dots, M-1\}$. Given conditional choice probabilities \mathbb{P} , I define the **representative value** in any state ω to be:

$$R(\omega | \mathbb{P}) = \sum_{n \in \mathcal{N}} \Pr(n)^{\frac{\lambda_1}{\lambda_M}} (\Pr(n | \mathcal{A}_1(\omega)))^{\frac{\lambda_2 - \lambda_1}{\lambda_M}} \dots (\Pr(n | \cap_{i=1}^{M-1} \mathcal{A}_i(\omega)))^{\frac{\lambda_M - \lambda_{M-1}}{\lambda_M}} e^{\frac{v_n(\omega)}{\lambda_M}}.$$

In the case of Shannon Entropy (when $\lambda_1 = \dots = \lambda_M = \lambda$), this representative value is much simpler since the weight used for each option is just its unconditional probability of being selected.

I next need to assign a score to n in each state of the world. This is easier to do. I define the **score** of option n in state ω to be:

$$s_n(\omega | \mathbb{P}) = \frac{e^{\frac{v_n(\omega)}{\lambda_M}}}{R(\omega | \mathbb{P})}.$$

Corollary 1.

Conditional and unconditional choice probabilities described in [\(5\)](#) and [\(6\)](#) are a solution to

(7) subject to (8) and (9) iff they comply with [Theorem 2](#) and solve:

$$\max_{\mathbb{P}} \sum_{\omega \in \Omega} \log \left(\sum_{n \in \mathcal{N}} \Pr(n)^{\frac{\lambda_1}{\lambda_M}} \Pr(n|\mathcal{A}_1(\omega))^{\frac{\lambda_2 - \lambda_1}{\lambda_M}} \dots \Pr(n|\cap_{i=1}^{M-1} \mathcal{A}_i(\omega))^{\frac{\lambda_M - \lambda_{M-1}}{\lambda_M}} e^{\frac{v_n(\omega)}{\lambda_M}} \right) \mu(\omega),$$

such that:

$$\forall A \in \mathcal{F} : \Pr(n|A) \geq 0 \quad \forall n, \quad \text{and} \quad \sum_{n \in \mathcal{N}} \Pr(n|A) = 1.$$

[Corollary 1](#) is helpful because it reduces the number of choice variables faced by the agent, which means it is easier for the researcher to find optimal agent behavior.

Theorem 3.

Conditional and unconditional choice probabilities described in (5) and (6) are a solution to (7) subject to (8) and (9) iff they comply with [Theorem 2](#), and $\forall n \in \mathcal{N}$ and $\forall \tilde{\omega} \in \Omega$: if $n \in \mathcal{C}(\mathbb{P})$ then $\Pr(n|\tilde{\omega}) > 0$, and if $n \notin \mathcal{C}(\mathbb{P})$ then

$$\begin{aligned} & \frac{\lambda_1}{\lambda_M} \left(\sum_{\omega \in \Omega} s_n(\omega|\mathbb{P}) \mu(\omega) \right) + \frac{\lambda_2 - \lambda_1}{\lambda_M} \left(\sum_{\omega \in \mathcal{A}_1(\tilde{\omega})} s_n(\omega|\mathbb{P}) \mu(\omega|\mathcal{A}_1(\omega)) \right) \\ & + \dots + \frac{\lambda_M - \lambda_{M-1}}{\lambda_M} \left(\sum_{\omega \in \cap_{i=1}^{M-1} \mathcal{A}_i(\tilde{\omega})} s_n(\omega|\mathbb{P}) \mu(\omega|\cap_{i=1}^{M-1} \mathcal{A}_i(\omega)) \right) \leq 1. \end{aligned}$$

[Theorem 3](#) establishes the MASE analogue of [Caplin et al. \(2018\)](#)'s Proposition 1, their central proposition. [Theorem 3](#) thus establishes necessary and sufficient conditions for optimal behavior in settings where MASE is used to measure the cost of information.

As is true with standard Shannon Entropy, optimal choice behavior may not be unique. If two options are known *a priori* to take the same value in each state of the world, for instance, then the agent can shift probability from one of these two options to the other whenever the former has a strictly positive probability of being selected in an optimal solution. While these sorts of environments are possible, optimal behavior is unique generically. This feature of optimal behavior should be evident since payoffs are linear, and costs are strictly convex. The exact sufficient conditions for the uniqueness of a solution are withheld, but for the solution not to be unique, similar to the case with Shannon Entropy studied by [Matějka and McKay \(2015\)](#), a very rigid form of co-movement is required between payoffs and states.

5 Comparisons with the Standard Model

In this section I compare and contrast the choice behavior that is produced by RI with Shannon Entropy and the choice behavior produced by the RI model developed in [Section 4](#) that uses the MASE measure developed in [Section 3](#). I first discuss the relationship between RU models and RI with MASE, and then revisit the two motivating examples, [Example 1](#) and [Example 2](#).

5.1 Comparison with Random Utility Model

It is standard practice to use a RU model to describe discrete choice settings. In such a model, the agent picks the option with the largest sum $u_n = v_n + \epsilon_n$ over all options $n \in \mathcal{N}$. Generally, u_n represents the value of the option to the agent, v_n represents the average value of the option across agents, and ϵ_n represents an idiosyncratic value to the agent. The role ϵ_n plays is up to interpretation, however, and is determined by the researchers specification ([Train, 2009](#)). In a setting where agents are thought to be rationally inattentive, the above terms are interpreted in a different way because the agent’s noisy behavior is generated by perceptual error instead of idiosyncratic differences in taste. In such settings, u_n represents the perceived value to the agent, v_n represents the true value to the agent, and ϵ_n is interpreted as an unobservable perceptual error that results from the noisy information strategy selected by the agent. [Woodford \(2014\)](#) argues that this latter interpretation is necessary in many contexts due to the stochastic responses observed in perceptual discrimination tasks such as those administered by [Dean and Neligh \(2018\)](#), which are akin to [Example 1](#) in [Section 2.1](#). While the interpretation of ϵ_n is relevant for welfare analysis, it is inconsequential for the description of choice behavior. How then can MASE be interpreted in terms of an RU framework, and what insights may be provided about the fitting of RU models?

[Matějka and McKay \(2015\)](#) point out that choice probabilities predicted by RI with Shannon Entropy correspond to multinomial logit choice probabilities where it is as if option values have been shifted due to the agent’s prior about potential values. An option that seems more desirable *a priori* is more likely to be selected by the agent in every state of the world, and thus is overvalued by a multinomial logit regression.

Rational inattention with MASE takes this one step further, as is shown by [Theorem 4](#), allowing the shift in perceived value to also depend on easier to observe attributes (attributes associated multipliers that are less than λ_M). This flexibility seems natural in many real world environments. Consider an agent that is trying to select a restaurant to go to. One may expect

that the chance of the agent selecting a given option to increase not only with the quality of the restaurant, and their prior impression of it, but also with easy to observe information such as on-line ratings the restaurant may have received.

Theorem 4:

The choice behavior described by \mathbb{P} , a solution to (7) subject to (8) and (9), is identical to the behavior produced by an RU model where each option $n \in \mathcal{N}$ has perceived value:

$$u_n = \tilde{v}_n + \alpha_n + \epsilon_n,$$

where $\tilde{v}_n = \frac{\mathbf{v}_n(\omega)}{\lambda_M}$, ϵ_n has an iid Gumbel distribution, and:

$$\alpha_n = \frac{\lambda_1}{\lambda_M} \log(N\Pr(n)) + \frac{\lambda_2 - \lambda_1}{\lambda_M} \log(N\Pr(n|\mathcal{A}_1(\omega))) + \dots + \frac{\lambda_M - \lambda_{M-1}}{\lambda_M} \log(N\Pr(n|\cap_{i=1}^{M-1} \mathcal{A}_i(\omega))).$$

[Theorem 4](#) is meant to provide insight into the outcome of attempting to fit a RU model in an environment where agents are rationally inattentive with a cost function for information described by MASE. [Theorem 4](#) does not say that a model of RI with MASE is equivalent to a RU model. Even if choice data from a given choice problem cannot be used to reject one for the other, across choice problems MASE produces behavior that can reject the hypothesis of a RU model. With MASE, for instance, as with standard Shannon Entropy, adding an option can increase the chance of an existing option being selected, which is not possible with a RU model.

Also, it is worth mentioning that since optimal behavior may result in some options being selected with probability zero, [Theorem 4](#) implicitly defines each α_n on the extended reals so that $\alpha_n = -\infty$ if $\Pr(n) = 0$.¹⁸

5.2 Example 1 Revisited

I now revisit [Example 1](#) from [Section 2.1](#), which is described in [Table 1](#). It seems natural that it should be easier for the agent to answer the question ‘Are 60 of the balls blue?’, than it is for them to answer ‘Are 51 or more of the balls blue?’. Similarly, it seems natural that it should be easier for the agent to answer the question ‘Are 60 of the balls red?’, than it is for them to answer ‘Are 51 or more of the balls red?’. Symmetry also means that the questions ‘Are 60 of the balls blue?’ and ‘Are 60 of the balls red?’ should have the same expected cost, and the questions ‘Are 51 or more

¹⁸It is shown in the proof of [Theorem 3](#) that if optimal behavior results in $\Pr(n) > 0$, then $\Pr(n|\omega) > 0 \forall \omega \in \Omega$.

of the balls blue?’ and ‘Are 51 or more of the balls red?’ should have the same expected cost. I can thus assume $\mathcal{A}_1 = \{A_1, A_2, A_3\} = \{\{\omega_1\}, \{\omega_2 \cup \omega_3\}, \{\omega_4\}\}$, and $\mathcal{A}_2 = \{\{\omega_1 \cup \omega_2\}, \{\omega_3 \cup \omega_4\}\}$.

Solutions to [Corollary 1](#) combined with [Theorem 2](#) mean that the chance of the agent selecting option 1 is increasing in the number of blue balls, as can be seen in [Figure 1](#), which depicts optimal behavior in each state of the world for a range of λ_1 . When λ_1 is small relative to λ_2 the agent chooses option 1 in state ω_1 with a high probability, and choose option 2 in state ω_4 with a high probability. The agent is thus better able to discern the state of the world when there are 40 of one color ball and 60 of the other than when there are 49 of one color and 51 of the other. This is supported by the experimental work of [Dean and Neligh \(2018\)](#), and is in contrast with the behavior predicted by a model of RI that uses Shannon Entropy.

[Morris and Yang \(2016\)](#) identify a related issue with Shannon Entropy’s lack of perceptual distance, and warn against its use in some continuous settings because it predicts discontinuous changes in behavior at places where payoffs change discontinuously. In the limit, as the number of different perceptual distances is allowed to grow, MASE can be used to produce the kind of continuous behavior that [Morris and Yang \(2016\)](#) desire.

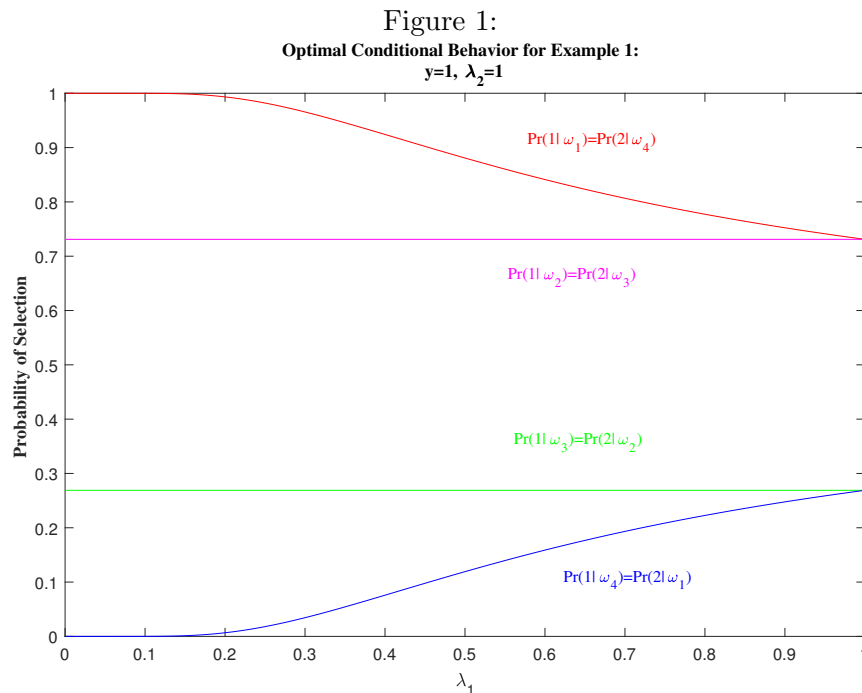
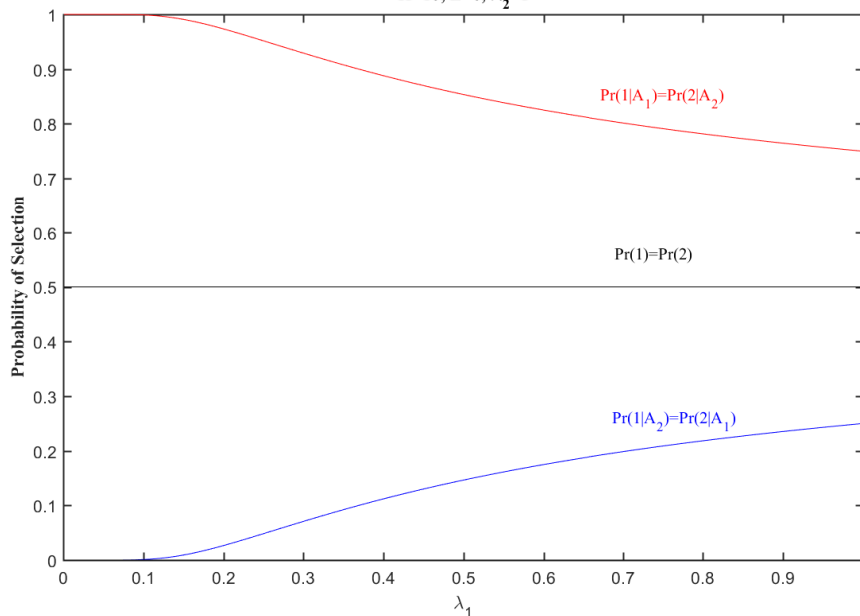


Figure 2:
Solutions to Corollary 1 for Example 2:
 $H=10, L=0, \lambda_2=1$



5.3 Example 2 Revisited

I now revisit [Example 2](#) from [Section 2.2](#), which is described in [Table 2](#). I assumed that learning the value of option 1 is less costly than learning the value of option 2. That is to say there are two attributes of the choice environment, one determines the value of option 1, the other determines the value of option 2, and the attribute that determines the value of option 1 is less costly to learn about. I can thus assume: $\mathcal{A}_1 = \{A_1, A_2\} = \{\{\omega_1 \cup \omega_2\}, \{\omega_3 \cup \omega_4\}\}$, and $\mathcal{A}_2 = \{\{\omega_1 \cup \omega_3\}, \{\omega_2 \cup \omega_4\}\}$.

Solutions to [Corollary 1](#) in this environment for a range of λ_1 can be found in [Figure 2](#), which shows that when λ_1 is small compared to λ_2 , the agent selects option 1 with a high probability when it is of value H , and selects option 2 with a high probability when option 1 is of value L . As λ_1 increases relative to λ_2 , the chance of option 1 being selected when it is of value H decreases. Similarly, as λ_1 increases relative to λ_2 , the chance of option 1 being selected when it is of value L increases. Note that the solutions to [Corollary 1](#) mean that the agent is more likely to select option 1 when state ω_1 has been realized since $\Pr(1|A_1) > \Pr(2|A_1)$, and more likely to select option 2 when state ω_4 has been realized since $\Pr(1|A_2) < \Pr(2|A_2)$, as can be observed with [Theorem 2](#).

Solutions to [Corollary 1](#) combined with [Theorem 4](#) mean that if an econometrician tries to fit this environment with a multinomial logit model that their estimate of H_1 , the high value of option

1, is biased upwards by $\frac{\lambda_2 - \lambda_1}{\lambda_2} \log(2\Pr(1|A_1))$, which is greater than zero since $\Pr(1|A_1) > 1/2$, and their estimate of L_1 , the low value of option 1, is biased downwards by $\frac{\lambda_2 - \lambda_1}{\lambda_2} \log(2\Pr(1|A_2))$, which is less than zero since $\Pr(1|A_2) < 1/2$. These biases are despite the fact that the unconditional chance of either option being selected is the same: $\Pr(1) = \Pr(2) = 1/2$. As such, the econometrician may have believed their analysis was not susceptible to informational biases if they had used Shannon Entropy to model the environment.

Further, as was mentioned, RU models and RI models with Shannon Entropy can both be rejected for RI with MASE in this environment if it is possible to alter the correlation between the values of the two options while holding the marginal distributions over values fixed for each option.¹⁹ If a RU model describes the agent, then changing the correlation between the values of the two options would not change the choice behavior of the agent in any state. If the behavior of the agent is instead described by MASE, then changing the correlation between the values of the two options would change the choice behavior of the agent in individual states because the total information that can be acquired from learning the value of option 1 (the option that is easier to learn about) changes with the correlation of the options' values. Further, if the above MASE specification is correct, the unconditional choice probabilities of the agent would remain constant when correlation is changed due to the symmetry of the environment, as long as the agent is doing some learning.²⁰ Finally, if the behavior of the agent is instead described by Shannon Entropy, in contrast, then the choice behavior in the individual states could only change if the unconditional choice probabilities changed.

6 Literature Review

Shannon Entropy has been used in several contexts to demonstrate informational biases in RU models. [Matějka and McKay \(2015\)](#) use Shannon Entropy in a model of RI to demonstrate the potential for informational biases in the multinomial logit model, while [Steiner, Stewart, and Matějka \(2017\)](#) use Shannon Entropy in a model of RI to demonstrate the potential for a similar bias in dynamic Logit. These results are significant for those who wish to fit RU models because, while observational data may coincide with the assumptions of a fitted RU model, informational biases can potentially invalidate counterfactual and welfare analysis, two common goals of such a

¹⁹This assertion and the assertions that follow in this paragraph are not difficult to show with [Theorem 2](#) and [Corollary 1](#).

²⁰The agent is doing some learning if their choice probabilities differ at all in states of the world that are realized with positive probability.

fitting.

The Shannon RI model has also led to a number of predictive successes. [Acharya and Wee \(2020\)](#) show that using Shannon Entropy to model firms as rationally inattentive results in a better fitting of labor market dynamics after the great depression. [Dasgupta and Mondria \(2018\)](#) show that using Shannon Entropy to model importers as rationally inattentive results in novel predictions that are supported by trade data. [Ambuehl, Ockenfels, and Stewart \(2019\)](#) experimentally verify predictions of Shannon Entropy in environments where agents are rationally inattentive to the consequences of participating in different transactions.

Perhaps as a response to the success Shannon Entropy has enjoyed, several recent papers have noted that Shannon Entropy may be a poor measure of the cost of acquiring information in some environments ([Caplin et al., 2017](#); [Morris & Yang, 2016](#)) because it lacks what is called “perceptual distance” ([Caplin et al., 2017](#), p. 39). As was alluded to previously, these papers argue that (i) more similar outcomes (outcomes that have less perceptual distance between them) should be more difficult to differentiate between, and (ii) when this property is missing, predicted behavior can differ significantly from the type of behavior that it would seem natural to expect ([Morris & Yang, 2016](#); [Dean & Neligh, 2018](#)).

An ad hoc group of cost functions that generalize Shannon Entropy is provided by [Huettner et al. \(2019\)](#). In their paper, the different options that the agent can choose between are allowed to differ in how costly they are to learn about. However, the group of cost functions developed by [Huettner et al. \(2019\)](#) are a strict subset of the cost functions that can be defined with MASE and do not solve the problems identified by [Morris and Yang \(2016\)](#) and [Dean and Neligh \(2018\)](#) that are solved by MASE (see [Example 1](#) in [Section 2.1](#)).

To better understand the relationship between the cost of learning and agent behavior, a number of papers have studied axiomatic models of rational inattention. Different papers, however, choose to focus their axioms on different aspects of the choice environment. [Caplin et al. \(2017\)](#), for instance, develop axioms that focus on the choice behavior of an agent after they expend effort to learn about the state of the world. In contrast, [de Oliveira \(2014\)](#) and [de Oliveira, Denti, Mihm, and Ozbek \(2017\)](#) develop axioms that focus on an agent’s preferences over choice menus before they expend effort to learn about the state of the world. Broadly, these papers aim to understand what implications rational agent behavior has for the form of information cost functions.

[Ellis \(2018\)](#) features axioms that focus on choice behavior and studies the implications for information cost functions, but further assumes that the agent learns by picking a partition of the

state space. While MASE uses the cost of learning the realized event of partitions as a primitive, the model studied in this paper does not constrain agents so that they must learn using partitions of the state space, and it can be shown that in a model of RI with MASE it is never optimal for the agent to choose an information strategy that is equivalent to a partition of the state space.²¹

Closer in nature to the work done in this paper, [Pomatto et al. \(2019\)](#) develop axioms that focus directly on the costs of information. Axioms that focus on costs for information are interesting because intuitive properties for costs of information can lead to unintuitive agent behavior that is compelling given real-world observations ([Gigerenzer & Todd, 1999](#)), but is often mistaken for irrational when axioms that appear rational are imposed on behavior. MASE, for instance, predicts ‘non-compensatory’ behavior, whereby changing an option so that it is more valuable to the agent can result in a lower chance of it being selected. This type of behavior raises important questions for welfare and counterfactual analysis, making effective policy design more challenging.

Unlike the work of [Pomatto et al. \(2019\)](#), which features axioms that are concerned with probabilistic experiments that can result in different outcomes in the same state of the world, this paper’s axioms are concerned with deterministic experiments (questions) that always result in the same outcome in a given state of the world, and contradict the form of constant marginal cost assumed in their paper.

7 Conclusion

Models of rational inattention that use Shannon Entropy to measure the cost of learning can help to better fit observed data in a range of contexts and also demonstrate that informational biases in random utility models can be significant for welfare and counterfactual analysis. While Shannon Entropy is a flexible and tractable tool, it does not allow for the attributes of the options an agent is choosing between to differ in their associated learning costs, which limits its application in economic environments.

This paper contributes to the literature by proposing and axiomatizing a new measure of uncertainty, Multi-Attribute Shannon Entropy (MASE), which is described in [Theorem 1](#), that allows for the different attributes of the options faced by an agent to differ in their associated learning costs. Like Shannon’s original axioms, the axioms in this paper focus on simple questions, ones that can be represented by partitions of the state space the agent is learning about. This

²¹This is true whenever the agent does some costly learning.

means the axioms in this paper are relatively easy to understand. Further, MASE can still be used as a cost function for information in a general model of rational inattention where the agent can choose to learn through any signal structure they desire.

MASE is shown to be a natural multi-parameter generalization of Shannon Entropy that maintains much of the tractability of Shannon’s standard measure. [Theorem 2](#) establishes the MASE analogue of [Matějka and McKay \(2015\)](#)’s necessary conditions for optimal behavior in the context of Shannon Entropy, and [Theorem 3](#) establishes the MASE analogue of [Caplin et al. \(2018\)](#)’s necessary and sufficient conditions for optimal behavior in the context of Shannon Entropy.

MASE also identifies a new form of informational bias demonstrated in [Theorem 4](#). The new form of bias can be present even when the agent has the same probability of selecting each option, which may seem to indicate an unbiased environment based on the previous literature. The biases that have previously been identified in the literature are independent of the realized state of the world, depending only on the agent’s prior about the environment. The informational biases that MASE identify are caused by attributes varying in their associated learning costs and can result in the same option being overvalued by a multinomial logit model for some realizations of its attributes and undervalued for other realizations of its attributes.

References

- Acharya, S., & Wee, S. L. (2020). Rational inattention in hiring decisions. *American Economic Journal: Macroeconomics*, 12(1), 1–40.
- Ambuehl, S., Ockenfels, A., & Stewart, C. (2019). Attention and selection effects. *Rotman School of Management Working Paper*(3154197).
- Caplin, A., Dean, M., & Leahy, J. (2017). *Rationally inattentive behavior: Characterizing and generalizing shannon entropy* (Tech. Rep.). National Bureau of Economic Research.
- Caplin, A., Dean, M., & Leahy, J. (2018). Rational inattention, optimal consideration sets, and stochastic choice. *The Review of Economic Studies*, 86(3), 1061–1094.
- Dasgupta, K., & Mondria, J. (2018). Inattentive importers. *Journal of International Economics*, 112, 150–165.
- Dean, M., & Neligh, N. L. (2018). Experimental tests of rational inattention.
- de Oliveira, H. (2014). *Axiomatic foundations for entropic costs of attention* (Tech. Rep.). Mimeo.
- de Oliveira, H., Denti, T., Mihm, M., & Ozbek, K. (2017). Rationally inattentive preferences and hidden information costs. *Theoretical Economics*, 12(2), 621–654.
- Ellis, A. (2018). Foundations for optimal inattention. *Journal of Economic Theory*, 173, 56–94.
- Gigerenzer, G., & Todd, P. M. (1999). Fast and frugal heuristics: The adaptive toolbox. In *Simple heuristics that make us smart* (pp. 3–34). Oxford University Press.
- Hébert, B., & Woodford, M. (2017). *Rational inattention and sequential information sampling* (Tech. Rep.). National Bureau of Economic Research.
- Huettner, F., Boyacı, T., & Akçay, Y. (2019). Consumer choice under limited attention when alternatives have different information costs. *Operations Research*.
- Lange, K. (2013). *Optimization* (Second Edition ed.; G. Casella, I. Olkin, & S. Fienberg, Eds.). Springer.
- Mackowiak, B., Matejka, F., & Wiederholt, M. (2021). Rational inattention: A review.
- Matějka, F., & McKay, A. (2015). Rational inattention to discrete choices: A new foundation for the multinomial logit model. *American Economic Review*, 105(1), 272–98.
- Morris, S., & Strack, P. (2019). The wald problem and the relation of sequential sampling and ex-ante information costs.
- Morris, S., & Yang, M. (2016). Coordination and continuous choice. *Working paper*.
- Noguchi, T., & Stewart, N. (2014). In the attraction, compromise, and similarity effects, alternatives

- are repeatedly compared in pairs on single dimensions. *Cognition*, 132(1), 44–56.
- Noguchi, T., & Stewart, N. (2018). Multialternative decision by sampling: A model of decision making constrained by process data. *Psychological review*, 125(4), 512.
- Pomatto, L., Strack, P., & Tamuz, O. (2019). The cost of information.
- Shannon, C. E. (1948). A mathematical theory of communication. *The Bell System Technical Journal*, 27(3), 379–423.
- Sims, C. A. (2003). Implications of rational inattention. *Journal of monetary Economics*, 50(3), 665–690.
- Steiner, J., Stewart, C., & Matějka, F. (2017). Rational inattention dynamics: Inertia and delay in decision-making. *Econometrica*, 85(2), 521–553.
- Stewart, N., Chater, N., & Brown, G. D. (2006). Decision by sampling. *Cognitive psychology*, 53(1), 1–26.
- Train, K. E. (2009). *Discrete choice methods with simulation*. Cambridge university press.
- Woodford, M. (2014). Stochastic choice: An optimizing neuroeconomic model. *American Economic Review*, 104(5), 495–500.

Appendix 1

Before I prove [Lemma 1](#), I show some other useful results:

Lemma 7. If a partition $\tilde{\mathcal{P}}$ is coarser than a learning strategy invariant partition \mathcal{P} , then $\tilde{\mathcal{P}}$ is also learning strategy invariant.

Proof. Suppose \mathcal{P} is a learning strategy invariant partition, and $\tilde{\mathcal{P}}$ is coarser than \mathcal{P} . If $\tilde{\mathcal{P}} = \mathcal{P}$ I am done.

If $\tilde{\mathcal{P}} \neq \mathcal{P}$, then the definition of learning strategy invariance tells us that for any learning strategy $\tilde{S} = (\mathcal{P}_1, \dots, \mathcal{P}_n)$ such that $\sigma(\tilde{S}) = \sigma(\tilde{\mathcal{P}})$, and any μ :

$$C(\mathcal{P}, \mu) = C((\tilde{\mathcal{P}}, \mathcal{P}), \mu) = C(\tilde{\mathcal{P}}, \mu) + \mathbb{E}[C(\mathcal{P}, \mu(\cdot|\tilde{\mathcal{P}}(\omega)))],$$

and,

$$C(\mathcal{P}, \mu) = C(\tilde{S}, \mu) + \mathbb{E}[C(\mathcal{P}, \mu(\cdot|\cap_{i=1}^n \mathcal{P}_i(\omega)))] = C(\tilde{S}, \mu) + \mathbb{E}[C(\mathcal{P}, \mu(\cdot|\tilde{\mathcal{P}}(\omega)))].$$

Thus, $C(\tilde{\mathcal{P}}, \mu) = C(\tilde{S}, \mu)$ for all such \tilde{S} , and any μ , so $\tilde{\mathcal{P}}$ is also learning strategy invariant. ■

Lemma 8. If $\mathcal{P} = \{A_1, \dots, A_m\}$ is a learning strategy invariant partition with $m \geq 3$, and probability measure μ assigns a probability of one to an event $A_i \in \mathcal{P}$, then $C(\mathcal{P}, \mu) = 0$.

Proof. Suppose $\mathcal{P} = \{A_1, \dots, A_m\}$ is a learning strategy invariant partition with $m \geq 3$ and there is an $A_i \in \mathcal{P}$ such that $\mu(A_i) = 1$. It is without loss to further assume $i = 1$.

Let $\tilde{\mathcal{P}} = \{A_1, A_1^c\}$, $\hat{\mathcal{P}} = \{A_1 \cup A_2, A_3, \dots, A_m\}$, $S_1 = (\tilde{\mathcal{P}}, \hat{\mathcal{P}})$, and $S_2 = (\tilde{\mathcal{P}}, \hat{\mathcal{P}}, \mathcal{P})$. The definition of learning strategy invariance tells us $C(S_1, \mu) = C(S_2, \mu)$, so $C(\mathcal{P}, \mu) = 0$. ■

Proof of Lemma 1. Suppose $\mathcal{P} = \{A_1, \dots, A_m\}$ is a learning strategy invariant partition of the state space Ω with $m \geq 3$. The definition of learning strategy invariance implies $C(\mathcal{P}, \mu)$ is fully determined by expected learning costs on binary partitions coarser than \mathcal{P} . [Axiom 1](#) tells us knowing $\mu(A_1), \dots$, and $\mu(A_m)$ is more than enough to compute expected learning costs on binary partitions coarser than \mathcal{P} , and thus $C(\mathcal{P}, \mu)$ is fully determined by $\mu(A_1), \dots$, and $\mu(A_m)$.

If I then show that for any $i, j \in \{1, \dots, m\}$ such that $i \neq j$, $C(\mathcal{P}, \mu) = C(\mathcal{P}, \tilde{\mu})$ if $\mu(A_k) = \tilde{\mu}(A_k)$ for $k \notin \{i, j\}$, $\mu(A_i) = \tilde{\mu}(A_j)$, and $\mu(A_j) = \tilde{\mu}(A_i)$, then the desired result holds, since a series of pairwise switches like this can be used to create any permutation desired. Assume that μ and $\tilde{\mu}$ satisfy the conditions from the previous sentence. It is without loss to assume $i = 1$ and $j = 2$. Define $\tilde{\mathcal{P}} = \{A_1, A_2, (A_1 \cup A_2)^c\}$. Notice that $\tilde{\mathcal{P}}$ must be learning strategy invariant based on [Lemma 7](#). Further, if I show that $C(\tilde{\mathcal{P}}, \mu) = C(\tilde{\mathcal{P}}, \tilde{\mu})$, then $C(\mathcal{P}, \mu) = C(\mathcal{P}, \tilde{\mu})$, since, if I

define $\hat{\mathcal{P}} = \{A_1 \cup A_2, A_3, \dots, A_m\}$, which is also learning strategy invariant based on [Lemma 7](#), then [Lemma 8](#) tells us:

$$\begin{aligned} C(\mathcal{P}, \mu) &= C(\tilde{\mathcal{P}}, \mu) + (1 - \mu(A_1 \cup A_2))C(\hat{\mathcal{P}}, \hat{\mu}) \\ &= C(\tilde{\mathcal{P}}, \tilde{\mu}) + (1 - \mu(A_1 \cup A_2))C(\hat{\mathcal{P}}, \hat{\mu}) = C(\mathcal{P}, \tilde{\mu}), \end{aligned}$$

if I define probability measure $\hat{\mu}$ so that $\hat{\mu}(A_1) = \hat{\mu}(A_2) = 0$, and for $i \in \{3, \dots, m\}$ I have $\hat{\mu}(A_i) = \mu(A_i)/(1 - \mu(A_1 \cup A_2))$. Now, let $\mathcal{P}_1^b = \{A_1, A_1^c\}$, $\mathcal{P}_2^b = \{A_2, A_2^c\}$, and $\mathcal{P}_3^b = \{A_1 \cup A_2, (A_1 \cup A_2)^c\}$. Notice \mathcal{P}_1^b , \mathcal{P}_2^b and \mathcal{P}_3^b , are all coarser than $\tilde{\mathcal{P}}$. Then, since $\tilde{\mathcal{P}}$ is learning strategy invariant, [Lemma 7](#) and [Lemma 8](#) tell us:

$$C(\tilde{\mathcal{P}}, \mu) = C(\mathcal{P}_3^b, \mu) + \mathbb{E}[C(\mathcal{P}_1^b, \mu(\cdot|\mathcal{P}_3^b(\omega)))], \text{ and } C(\tilde{\mathcal{P}}, \tilde{\mu}) = C(\mathcal{P}_3^b, \tilde{\mu}) + \mathbb{E}[C(\mathcal{P}_1^b, \tilde{\mu}(\cdot|\mathcal{P}_3^b(\omega)))].$$

Notice that [Axiom 1](#) tells us $C(\mathcal{P}_3^b, \mu) = C(\mathcal{P}_3^b, \tilde{\mu})$ since both μ and $\tilde{\mu}$ assign the same probability to the events $A_1 \cup A_2$ and $(A_1 \cup A_2)^c$. So, all that remains to show is that if the probability measure $\tilde{\nu}$ is a permutation of the probability measure ν on \mathcal{P}_1^b , then $C(\mathcal{P}_1^b, \nu) = C(\mathcal{P}_1^b, \tilde{\nu})$. Fix arbitrary $\nu(A_1) = x \in [0, 1]$. Now consider the probability measures q_1, q_2, q_3 , such that:

$$q_1(A_1) = x, \quad q_1(A_2) = 0, \quad q_1((A_1 \cup A_2)^c) = 1 - x,$$

$$q_2(A_1) = 0, \quad q_2(A_2) = x, \quad q_2((A_1 \cup A_2)^c) = 1 - x,$$

$$q_3(A_1) = 1 - x, \quad q_3(A_2) = x, \quad q_3((A_1 \cup A_2)^c) = 0.$$

Notice that q_3 is a permutation of q_1 on \mathcal{P}_1^b . So then, using [Axiom 1](#), the definition of learning strategy invariance, and [Lemma 8](#), all repeatedly:

$$\begin{aligned} C(\mathcal{P}_1^b, q_1) &= C(\tilde{\mathcal{P}}, q_1) = C(\mathcal{P}_3^b, q_1) = C(\mathcal{P}_3^b, q_2) \\ &= C(\tilde{\mathcal{P}}, q_2) = C(\mathcal{P}_2^b, q_2) = C(\mathcal{P}_2^b, q_3) = C(\tilde{\mathcal{P}}, q_3) = C(\mathcal{P}_1^b, q_3), \end{aligned}$$

and I am done. ■

Proof of [Lemma 2](#). For all partitions $\mathcal{P} = \{A_1, \dots, A_m\}$ and probability measures μ defined on \mathcal{P} , define the vector $\mu(\mathcal{P}) = (\mu(A_1), \dots, \mu(A_m))$.

Suppose C satisfies [Axiom 1](#), that $\mathcal{P}_i = \{A_1, \dots, A_m\}$ is a learning strategy invariant with $m \geq 3$, and $\tilde{\mathcal{P}}_i$ is another learning strategy invariant partition that is coarser than \mathcal{P}_i . [Lemma 1](#) tells us that $C(\mathcal{P}_i, \mu)$ is fully determined by $\mu(\mathcal{P}_i)$, and if the strictly positive entries of $\mu(\mathcal{P}_i)$ and $\mu(\tilde{\mathcal{P}}_i)$ are the same (up to a permutation), then the addition of [Lemma 8](#) and the definition of learning strategy invariant partitions tell us $C(\mathcal{P}_i, \mu) = C(\tilde{\mathcal{P}}_i, \mu)$ since I can pick μ so that uncertainty about which event in \mathcal{P}_i has been realized is fully determined by the realized event of $\tilde{\mathcal{P}}_i$. What does this mean? This means that there is a function which maps from vectors of probabilities onto the reals, $c_i : \cup_{j=1}^{m-1} \Delta^j \rightarrow \mathbb{R}$, where Δ^j is the j simplex, such that for any learning strategy invariant partition $\tilde{\mathcal{P}}_i$ coarser than \mathcal{P}_i , if the strictly positive entries of $\mu(\mathcal{P}_i)$ and $\mu(\tilde{\mathcal{P}}_i)$ are the same (up to a permutation) then $C(\tilde{\mathcal{P}}_i, \mu) = c_i(\mu(\tilde{\mathcal{P}}_i)) \equiv C(\mathcal{P}_i, \mu)$.

So, for any binary partition \mathcal{P}^b coarser than \mathcal{P}_i , $C(\mathcal{P}^b, \mu) = c_i(\mu(\mathcal{P}^b))$ (notice that this means that $C(\mathcal{P}^b, \mu)$ is constant with respect to permutations of μ on \mathcal{P}^b for all such \mathcal{P}^b since $C(\mathcal{P}, \mu)$ is constant with respect to permutations of μ on \mathcal{P}). Now pick $\tilde{\mathcal{P}}_i = \{B_1, B_2, B_3\}$ so that it is coarser than \mathcal{P}_i and it has three elements. [Lemma 7](#) tells us $\tilde{\mathcal{P}}_i$ is learning strategy invariant, and it is easy to show each binary partition which is coarser than $\tilde{\mathcal{P}}_i$ is coarser than \mathcal{P}_i . Thus, for all probability measures μ on $\tilde{\mathcal{P}}_i$ such that $\mu(B_1)$, $\mu(B_2)$, and $\mu(B_3)$ are all strictly less than one, the definition of learning strategy invariance tells us:

$$\begin{aligned} C(\tilde{\mathcal{P}}_i, \mu) &= c_i(\mu(B_1), 1 - \mu(B_1)) + (1 - \mu(B_1))c_i\left(\frac{\mu(B_2)}{\mu(B_2) + \mu(B_3)}, \frac{\mu(B_3)}{\mu(B_2) + \mu(B_3)}\right) \\ &= c_i(\mu(B_2), 1 - \mu(B_2)) + (1 - \mu(B_2))c_i\left(\frac{\mu(B_1)}{\mu(B_1) + \mu(B_3)}, \frac{\mu(B_3)}{\mu(B_1) + \mu(B_3)}\right) \\ &= c_i(\mu(B_3), 1 - \mu(B_3)) + (1 - \mu(B_3))c_i\left(\frac{\mu(B_1)}{\mu(B_1) + \mu(B_2)}, \frac{\mu(B_2)}{\mu(B_1) + \mu(B_2)}\right), \end{aligned}$$

and I am done. ■

I say that the vector (q_1, \dots, q_n) is a **permutation** of the vector (p_1, \dots, p_n) if there is a bijection $\pi : \{1, \dots, n\} \rightarrow \{1, \dots, n\}$ such that $\forall i \in \{1, \dots, n\}, q_i = p_{\pi(i)}$. Before I prove [Lemma 3](#), I pause to show two more useful results.

Lemma 9. Given a binary partition \mathcal{P}^b , if I define $c_{\mathcal{P}^b} : \cup_{j=1}^{\infty} \Delta^j \rightarrow \mathbb{R}$, where Δ^j is the j simplex, such that (for $n \geq 2$): $c_{\mathcal{P}^b}(p_1, \dots, p_n) = C(\mathcal{P}^b, p_1, 1 - p_1)$ if $p_1 + p_2 = 1$, and otherwise:

$$c_{\mathcal{P}^b}(p_1, \dots, p_n) = C(\mathcal{P}^b, p_1, 1 - p_1) + (1 - p_1)C\left(\mathcal{P}^b, \frac{p_2}{1 - p_1}, \frac{1 - p_1 - p_2}{1 - p_1}\right)$$

$$+ \dots + (1 - p_1 - \dots - p_{m-1})C\left(\mathcal{P}^b, \frac{p_m}{1 - p_1 - \dots - p_{m-1}}, \frac{1 - p_1 - \dots - p_m}{1 - p_1 - \dots - p_{m-1}}\right),$$

where m is the lowest integer such that $p_1 + \dots + p_m = 1$, then if (q_1, \dots, q_n) is a permutation of (p_1, \dots, p_n) , and C satisfies [Axiom 1](#), and [Axiom 2](#), then: $c_{\mathcal{P}^b}(q_1, \dots, q_n) = c_{\mathcal{P}^b}(p_1, \dots, p_n)$, and further if (p_1, \dots, p_n) is a vector ($n \geq 2$) with one entry of value one, and the rest zero $c_{\mathcal{P}^b}(p_1, \dots, p_n) = 0$.

Proof of Lemma 9. Given a binary partition \mathcal{P}^b , suppose C satisfies [Axiom 1](#), and [Axiom 2](#), and that $c_{\mathcal{P}^b}$ is defined as above. All vectors discussed in this proof are assumed to sum to one. I proceed with an inductive argument, beginning by showing $c_{\mathcal{P}^b}(p, 1 - p)$ satisfies the desired properties. Consider $c_{\mathcal{P}^b}(p_1, p_2, p_3)$ when $p_1, p_3 > 0$, and $p_2 = 0$. [Axiom 2](#) tells us:

$$c_{\mathcal{P}^b}(p_1, 1 - p_1) + (1 - p_1)c_{\mathcal{P}^b}(0, 1) = c_{\mathcal{P}^b}(0, 1) + c_{\mathcal{P}^b}(p_1, 1 - p_1) = c_{\mathcal{P}^b}(p_3, 1 - p_3) + (1 - p_3)c_{\mathcal{P}^b}(1, 0).$$

The first equality implies $c_{\mathcal{P}^b}(0, 1) = 0$. Now consider $c_{\mathcal{P}^b}(q_1, q_2, q_3)$ when $q_1, q_2 > 0$, and $q_3 = 0$. [Axiom 2](#) tells us:

$$c_{\mathcal{P}^b}(q_1, q_2) + (1 - q_1)c_{\mathcal{P}^b}(1, 0) = c_{\mathcal{P}^b}(0, 1) + c_{\mathcal{P}^b}(q_1, q_2),$$

so since $c_{\mathcal{P}^b}(0, 1) = 0$, I know $c_{\mathcal{P}^b}(1, 0) = 0 = c_{\mathcal{P}^b}(0, 1)$, and combined with the previous two equalities above I know:

$$c_{\mathcal{P}^b}(p_1, 1 - p_1) = c_{\mathcal{P}^b}(p_3, 1 - p_3) + (1 - p_3)c_{\mathcal{P}^b}(1, 0) = c_{\mathcal{P}^b}(1 - p_1, p_1).$$

Thus, $c_{\mathcal{P}^b}(p, 1 - p) = c_{\mathcal{P}^b}(1 - p, p)$ for all $p \in [0, 1]$. Since $c_{\mathcal{P}^b}(1, 0) = 0$, when I show $c_{\mathcal{P}^b}$ is constant with respect to permutations of vectors of arbitrary length (greater or equal to two), it establishes that if (p_1, \dots, p_n) is a vector ($n \geq 2$) with one entry of value one, and the rest zero, then $c_{\mathcal{P}^b}(p_1, \dots, p_n) = 0$.

Next I show $c_{\mathcal{P}^b}(p_1, p_2, p_3)$ is constant with respect to permutations. Since $c_{\mathcal{P}^b}$ is constant with respect to permutation on vectors of length two, the definition of $c_{\mathcal{P}^b}$, and the fact that $c_{\mathcal{P}^b}(1, 0) = c_{\mathcal{P}^b}(0, 1) = 0$, tells us $c_{\mathcal{P}^b}(p_1, p_2, p_3) = c_{\mathcal{P}^b}(p_1, p_3, p_2)$. Thus, if I show for any probability vector of length three that $c_{\mathcal{P}^b}(p_1, p_2, p_3) = c_{\mathcal{P}^b}(p_2, p_1, p_3)$, then $c_{\mathcal{P}^b}(p_1, p_2, p_3)$ is constant with respect to permutations since combinations of these two different pairwise permutations can achieve any permutation desired. This is easy to show since if $p_1 = 1$ or $p_2 = 1$, then I know this

is true, and otherwise with [Axiom 2](#) I know:

$$\begin{aligned} c_{\mathcal{P}^b}(p_1, p_2, p_3) &= c_{\mathcal{P}^b}(p_1, 1 - p_1) + (1 - p_1)c_{\mathcal{P}^b}\left(\frac{p_2}{1 - p_1}, \frac{1 - p_1 - p_2}{1 - p_1}\right) \\ &= c_{\mathcal{P}^b}(p_2, 1 - p_2) + (1 - p_2)c_{\mathcal{P}^b}\left(\frac{p_1}{1 - p_2}, \frac{1 - p_1 - p_2}{1 - p_2}\right) = c_{\mathcal{P}^b}(p_2, p_1, p_3). \end{aligned}$$

Now assume that $c_{\mathcal{P}^b}$ is constant with respect to permutations on vectors of length $n \geq 3$, and I next show $c_{\mathcal{P}^b}$ is constant with respect to permutations on vectors of length $n + 1$, and the proof is finished. If $p_1 + p_2 = 1$, then I am done. If not, notice that $c_{\mathcal{P}^b}(p_1, \dots, p_{n+1}) = c_{\mathcal{P}^b}(p_1, 1 - p_1) + (1 - p_1)c_{\mathcal{P}^b}\left(\frac{p_2}{1 - p_1}, \dots, \frac{p_{n+1}}{1 - p_1}\right)$, whenever $p_1 \neq 1$, and as part of the inductive argument I assumed $c_{\mathcal{P}^b}$ was constant with respect to permutations on vectors of length n , so I only need to show $c_{\mathcal{P}^b}(p_1, p_2, \dots, p_{n+1}) = c_{\mathcal{P}^b}(p_2, p_1, \dots, p_{n+1})$, which is true:

$$\begin{aligned} c_{\mathcal{P}^b}(p_1, p_2, \dots, p_{n+1}) &= c_{\mathcal{P}^b}(p_1, 1 - p_1) + (1 - p_1)c_{\mathcal{P}^b}\left(\frac{p_2}{1 - p_1}, \dots, \frac{p_{n+1}}{1 - p_1}\right) \\ &= c_{\mathcal{P}^b}(p_1, 1 - p_1) + (1 - p_1)c_{\mathcal{P}^b}\left(\frac{p_2}{1 - p_1}, \frac{1 - p_1 - p_2}{1 - p_1}\right) + (1 - p_1 - p_2)c_{\mathcal{P}^b}\left(\frac{p_3}{1 - p_1 - p_2}, \dots, \frac{p_{n+1}}{1 - p_1 - p_2}\right) \\ &= c_{\mathcal{P}^b}(p_1, p_2, 1 - p_1 - p_2) + (1 - p_1 - p_2)c_{\mathcal{P}^b}\left(\frac{p_3}{1 - p_1 - p_2}, \dots, \frac{p_{n+1}}{1 - p_1 - p_2}\right) \\ &= c_{\mathcal{P}^b}(p_2, p_1, 1 - p_1 - p_2) + (1 - p_1 - p_2)c_{\mathcal{P}^b}\left(\frac{p_3}{1 - p_1 - p_2}, \dots, \frac{p_{n+1}}{1 - p_1 - p_2}\right) \\ &= c_{\mathcal{P}^b}(p_2, 1 - p_2) + (1 - p_2)c_{\mathcal{P}^b}\left(\frac{p_1}{1 - p_2}, \frac{1 - p_1 - p_2}{1 - p_2}\right) + (1 - p_1 - p_2)c_{\mathcal{P}^b}\left(\frac{p_3}{1 - p_1 - p_2}, \dots, \frac{p_{n+1}}{1 - p_1 - p_2}\right) \\ &= c_{\mathcal{P}^b}(p_2, 1 - p_2) + (1 - p_2)c_{\mathcal{P}^b}\left(\frac{p_1}{1 - p_2}, \dots, \frac{p_{n+1}}{1 - p_2}\right) = c_{\mathcal{P}^b}(p_2, p_1, \dots, p_{n+1}). \blacksquare \end{aligned}$$

Lemma 10. Given a binary partition \mathcal{P}^b , define $c_{\mathcal{P}^b} : \cup_{j=1}^{\infty} \Delta^j \rightarrow \mathbb{R}$, where Δ^j is the j simplex, as in the statement of [Lemma 9](#), and suppose C satisfies [Axiom 1](#), and [Axiom 2](#), then if (q_1, \dots, q_m) and (p_1, \dots, p_n) are two probability vectors (weakly positive numbers that sum to one with $1 < m < n$), such that each q_i is strictly positive, and can be written as the sum of one or more p_j with each p_j used once in the sum of only one q_i . Let us rename the p_j (s) assigned to each q_i so that $q_i = p_1^i + \dots + p_{n_i}^i$. Then it is true that:

$$c_{\mathcal{P}^b}(p_1, \dots, p_n) = c_{\mathcal{P}^b}(q_1, \dots, q_m) + \sum_{i=1}^m q_i c_{\mathcal{P}^b}\left(\frac{p_1^i}{q_i}, \dots, \frac{p_{n_i}^i}{q_i}, 0\right).$$

Proof of Lemma 10. Given a binary partition \mathcal{P}^b , suppose C satisfies [Axiom 1](#), and [Axiom 2](#), that

$c_{\mathcal{P}^b}$ is defined as in the statement of [Lemma 9](#), and (q_1, \dots, q_m) and (p_1, \dots, p_n) are described as in the statement of [Lemma 10](#) (including the renaming of each p_j). All vectors discussed in this proof are assumed to sum to one and have at least length two, and I use the fact that the definition of $c_{\mathcal{P}^b}$ implies $c_{\mathcal{P}^b}(p_1, \dots, p_n) = c_{\mathcal{P}^b}(p_1, \dots, p_n, 0)$, and $c_{\mathcal{P}^b}(1, 0) = 0$, without reference. In [Lemma 9](#) I showed $c_{\mathcal{P}^b}$ is constant with respect to permutations of vectors of arbitrary length (greater or equal to two). Thus, all I need to do is show:

$$c_{\mathcal{P}^b}(p_1, \dots, p_{m-1}, p_m, \dots, p_n) = c_{\mathcal{P}^b}(q_1, \dots, q_m) + q_m c_{\mathcal{P}^b}\left(\frac{p_m}{q_m}, \dots, \frac{p_n}{q_m}, 0\right),$$

where for $i \in \{1, \dots, m-1\}$ $q_i = p_i$, $1 < m < n$, and $q_m = p_m + \dots + p_n > 0$. This is of course true. If $m = 2$, or $q_m = p_m$, this is trivially true. If $m > 2$ and $q_m > p_m$, then it is still true given the definition of $c_{\mathcal{P}^b}$ since (assuming without loss that $p_n > 0$):

$$\begin{aligned} c_{\mathcal{P}^b}(p_1, \dots, p_{m-1}, p_m, \dots, p_n) &= C(\mathcal{P}^b, p_1, 1 - p_1) + (1 - p_1)C\left(\mathcal{P}^b, \frac{p_2}{1 - p_1}, \frac{1 - p_1 - p_2}{1 - p_1}\right) \\ &+ \dots + (1 - p_1 - \dots - p_{m-1})C\left(\mathcal{P}^b, \frac{p_m}{1 - p_1 - \dots - p_{m-1}}, \frac{1 - p_1 - \dots - p_m}{1 - p_1 - \dots - p_{m-1}}\right) \\ &+ (1 - p_1 - \dots - p_m)C\left(\mathcal{P}^b, \frac{p_{m+1}}{1 - p_1 - \dots - p_m}, \frac{1 - p_1 - \dots - p_m}{1 - p_1 - \dots - p_m}\right) \\ &+ \dots + (1 - p_1 - \dots - p_{n-1})C\left(\mathcal{P}^b, \frac{p_n}{1 - p_1 - \dots - p_{n-1}}, \frac{1 - p_1 - \dots - p_n}{1 - p_1 - \dots - p_{n-1}}\right) \\ &= c_{\mathcal{P}^b}(q_1, \dots, q_m) + q_m c_{\mathcal{P}^b}\left(\frac{p_m}{q_m}, \dots, \frac{p_n}{q_m}, 0\right). \blacksquare \end{aligned}$$

Proof of [Lemma 3](#).

Given a binary partition $\mathcal{P}^b = \{A_1, A_2\}$, define $c_{\mathcal{P}^b} : \cup_{j=1}^{\infty} \Delta^j \rightarrow \mathbb{R}$, where Δ^j is the j simplex, as in the statement of [Lemma 9](#), and suppose C satisfies [Axiom 1](#), [Axiom 2](#), and [Axiom 3](#). Remember $C(\mathcal{P}^b, \mu) = c_{\mathcal{P}^b}(\mu(A_1), \mu(A_2))$ for all probability measures μ so [Lemma 9](#) tells us $C(\mathcal{P}^b, p, 1 - p) = C(\mathcal{P}^b, 1 - p, p)$, for each $p \in [0, 1]$, and I thus only wish to show $c_{\mathcal{P}^b}(p, 1 - p)$ is continuous for $p \in [0, 1]$. I proceed with a proof by contradiction: Suppose not, and $c_{\mathcal{P}^b}(p, 1 - p)$ is discontinuous at some point $p = p_d \in [0, 1]$. Since $c_{\mathcal{P}^b}(p, 1 - p) = c_{\mathcal{P}^b}(1 - p, p)$, it is without loss to assume $p_d \in [0, \frac{1}{2}]$.

First, notice that if $c_{\mathcal{P}^b}(p, 1 - p)$ is continuous at $p = 0$ then it is continuous at $p = \frac{1}{2}$: this is because [Axiom 2](#) tells us that for small $\delta > 0$: $c_{\mathcal{P}^b}(\delta, \frac{1}{2} - \frac{\delta}{2}, \frac{1}{2} - \frac{\delta}{2}) = c_{\mathcal{P}^b}(\delta, 1 - \delta) + (1 - \delta)c_{\mathcal{P}^b}(1/2, 1/2) = c_{\mathcal{P}^b}(\frac{1}{2} - \frac{\delta}{2}, \frac{1}{2} + \frac{\delta}{2}) + (\frac{1}{2} + \frac{\delta}{2})c_{\mathcal{P}^b}(\frac{2\delta}{1+\delta}, \frac{1-\delta}{1+\delta})$. Since [Axiom 3](#) requires that there

is some $p_c \in [0, \frac{1}{2}]$ such that $c_{\mathcal{P}^b}(p, 1-p)$ is continuous at p_c , it is thus without loss to assume $c_{\mathcal{P}^b}(p, 1-p)$ is continuous at $p_c \in (0, \frac{1}{2}]$.

Second, notice that it is not possible that the only $p \in [0, \frac{1}{2}]$ at which $c_{\mathcal{P}^b}(p, 1-p)$ is discontinuous is $p = 0$, because again [Axiom 2](#) tells us that for small $\delta > 0$: $c_{\mathcal{P}^b}(\delta, \frac{1}{2} - \frac{\delta}{2}, \frac{1}{2} - \frac{\delta}{2}) = c_{\mathcal{P}^b}(\delta, 1-\delta) + (1-\delta)c_{\mathcal{P}^b}(1/2, 1/2) = c_{\mathcal{P}^b}(\frac{1}{2} - \frac{\delta}{2}, \frac{1}{2} + \frac{\delta}{2}) + (\frac{1}{2} + \frac{\delta}{2})c_{\mathcal{P}^b}(\frac{2\delta}{1+\delta}, \frac{1-\delta}{1+\delta})$, and either:

$$\limsup_{p \downarrow 0} c_{\mathcal{P}^b}(p, 1-p) = H < \infty \text{ (with } H > 0) \text{ or } \limsup_{p \downarrow 0} c_{\mathcal{P}^b}(p, 1-p) = \infty.$$

If the former is true, then I can pick arbitrarily small $\delta \in (0, \frac{1}{4})$ to ensure that $c_{\mathcal{P}^b}(\delta, 1-\delta)$ is arbitrarily close to H , $c_{\mathcal{P}^b}(\frac{2\delta}{1+\delta}, \frac{1-\delta}{1+\delta})$ is not more than H , and $|(1-\delta)c_{\mathcal{P}^b}(1/2, 1/2) - c_{\mathcal{P}^b}(\frac{1}{2} - \frac{\delta}{2}, \frac{1}{2} + \frac{\delta}{2})| < \frac{1}{8}H$, which creates a contradiction. If, instead, the later is true, then I can pick arbitrarily small $\delta \in (0, \frac{1}{4})$ so that $c_{\mathcal{P}^b}(\delta, 1-\delta) \geq c_{\mathcal{P}^b}(p, 1-p) \forall p \in [\delta, \frac{1}{2}]$, and so that $|(1-\delta)c_{\mathcal{P}^b}(1/2, 1/2) - c_{\mathcal{P}^b}(\frac{1}{2} - \frac{\delta}{2}, \frac{1}{2} + \frac{\delta}{2})| < \frac{1}{8}c_{\mathcal{P}^b}(\delta, 1-\delta)$, which again creates a contradiction (remember that $\delta < \frac{2\delta}{1+\delta}$).

Third, if $c_{\mathcal{P}^b}(p, 1-p)$ is discontinuous at $p = \frac{1}{2}$ then it is discontinuous at a $p \in \{\frac{1}{4}, \frac{1}{3}\}$ because [Axiom 2](#) tells us that for small δ : $c_{\mathcal{P}^b}(\frac{1}{2} - \delta, \frac{1}{3} + \frac{2\delta}{3}, \frac{1}{6} + \frac{\delta}{3}) = c_{\mathcal{P}^b}(\frac{1}{2} - \delta, \frac{1}{2} + \delta) + (\frac{1}{2} + \delta)c_{\mathcal{P}^b}(\frac{1}{3}, \frac{2}{3}) = c_{\mathcal{P}^b}(\frac{1}{3} + \frac{2\delta}{3}, \frac{2}{3} - \frac{2\delta}{3}) + (\frac{2}{3} - \frac{2\delta}{3})c_{\mathcal{P}^b}((\frac{1}{6} + \frac{\delta}{3})/(\frac{2}{3} - \frac{2\delta}{3}), (\frac{1}{2} - \delta)/(\frac{2}{3} - \frac{2\delta}{3}))$. Thus it is without loss to assume $c_{\mathcal{P}^b}(p, 1-p)$ is discontinuous at $p_d \in (0, \frac{1}{2})$ (given second and third point).

It is not possible, however for $c_{\mathcal{P}^b}(p, 1-p)$ to be continuous at $p_c \in (0, \frac{1}{2}]$ and discontinuous at $p_d \in (0, \frac{1}{2})$. I can reach a contradiction in the following way: pick (p_1, p_2, p_3, p_4) such that they sum to one and:

$$p_1 + p_2 = p_d, \quad \frac{p_1}{p_1 + p_2} = p_c, \quad \text{and} \quad \frac{p_4}{p_3 + p_4} = p_c,$$

$$\text{so that as a result } p_1 + p_4 = p_c, \quad \frac{p_1}{p_1 + p_4} = p_d, \quad \text{and} \quad \frac{p_2}{p_2 + p_3} = p_d.$$

How these four probabilities are selected is quite important, and this is where a lot of the magic happens. Now, notice [Lemma 10](#) tells us:

$$\begin{aligned} & c_{\mathcal{P}^b}(p_1, p_2, p_3, p_4) \\ &= c_{\mathcal{P}^b}(p_1 + p_2, p_3 + p_4) + (p_1 + p_2)c_{\mathcal{P}^b}\left(\frac{p_1}{p_1 + p_2}, \frac{p_2}{p_1 + p_2}\right) + (p_3 + p_4)c_{\mathcal{P}^b}\left(\frac{p_3}{p_3 + p_4}, \frac{p_4}{p_3 + p_4}\right) \\ &= c_{\mathcal{P}^b}(p_1 + p_4, p_2 + p_3) + (p_1 + p_4)c_{\mathcal{P}^b}\left(\frac{p_1}{p_1 + p_4}, \frac{p_4}{p_1 + p_4}\right) + (p_2 + p_3)c_{\mathcal{P}^b}\left(\frac{p_2}{p_2 + p_3}, \frac{p_3}{p_2 + p_3}\right). \end{aligned}$$

Substituting in terms using the definitions of the four probabilities it is then clear that:

$$\begin{aligned} & c_{\mathcal{P}^b}(p_d, 1 - p_d) + (p_d)c_{\mathcal{P}^b}(p_c, 1 - p_c) + (1 - p_d)c_{\mathcal{P}^b}(1 - p_c, p_c) \\ &= c_{\mathcal{P}^b}(p_c, 1 - p_c) + (p_c)c_{\mathcal{P}^b}(p_d, 1 - p_d) + (1 - p_c)c_{\mathcal{P}^b}(p_d, 1 - p_d). \end{aligned}$$

Next, I would like to point out that $c_{\mathcal{P}^b}$ is discontinuous from both sides at p_d since I can increase p_1 and p_3 by a small $\delta > 0$, and decrease p_2 and p_4 by the same δ , and as δ is taken to zero, continuity at p_c implies the change in $c_{\mathcal{P}^b}(p_1 + p_2, p_3 + p_4) + (p_1 + p_2)c_{\mathcal{P}^b}\left(\frac{p_1}{p_1 + p_2}, \frac{p_2}{p_1 + p_2}\right) + (p_3 + p_4)c_{\mathcal{P}^b}\left(\frac{p_3}{p_3 + p_4}, \frac{p_4}{p_3 + p_4}\right)$ goes to zero, so discontinuities at either side of p_d must offset each other so the change in $c_{\mathcal{P}^b}(p_1 + p_4, p_2 + p_3) + (p_1 + p_4)c_{\mathcal{P}^b}\left(\frac{p_1}{p_1 + p_4}, \frac{p_4}{p_1 + p_4}\right) + (p_2 + p_3)c_{\mathcal{P}^b}\left(\frac{p_2}{p_2 + p_3}, \frac{p_3}{p_2 + p_3}\right)$ goes to zero with δ .

Next, I show that at p_d it must be that $c_{\mathcal{P}^b}(p, 1 - p)$ drops:

$$\liminf_{p \downarrow p_d} c_{\mathcal{P}^b}(p, 1 - p) < c_{\mathcal{P}^b}(p_d, 1 - p_d).$$

I proceed with a proof by contradiction: Suppose not, and that:

$$\liminf_{p \downarrow p_d} c_{\mathcal{P}^b}(p, 1 - p) \geq c_{\mathcal{P}^b}(p_d, 1 - p_d).$$

There are several cases of interest. In case one:

$$\liminf_{p \downarrow p_d} c_{\mathcal{P}^b}(p, 1 - p) = L > c_{\mathcal{P}^b}(p_d, 1 - p_d).$$

Case one is not possible, however, since I can choose arbitrarily small $\delta > 0$ and add it to p_1 and subtract it from p_4 so that $c_{\mathcal{P}^b}\left(\frac{p_1}{p_1 + p_4}, \frac{p_4}{p_1 + p_4}\right)$ is arbitrarily close to L , while $c_{\mathcal{P}^b}(p_1 + p_2, p_3 + p_4)$ is more than L , and all other terms remain essentially constant, creating a contradiction. In case two,

$$\liminf_{p \downarrow p_d} c_{\mathcal{P}^b}(p, 1 - p) = c_{\mathcal{P}^b}(p_d, 1 - p_d) \text{ and } \limsup_{p \downarrow p_d} c_{\mathcal{P}^b}(p, 1 - p) = H > c_{\mathcal{P}^b}(p_d, 1 - p_d),$$

with $H < \infty$. Case two is also not possible, however, since I can choose arbitrarily small $\delta > 0$ and add it to p_1 and subtract it from p_4 so that $c_{\mathcal{P}^b}(p_1 + p_2, p_3 + p_4)$ is arbitrarily close to H , while $c_{\mathcal{P}^b}\left(\frac{p_1}{p_1 + p_4}, \frac{p_4}{p_1 + p_4}\right)$ is less than H , and all other terms remain essentially constant, creating

a contradiction. Case three, the final case, is the same as case 2 except $H = \infty$. Case three is also not possible, however, since I can choose arbitrarily small $\delta > 0$ and add it to p_1 and p_3 and subtract it from p_2 and p_4 so that $c_{\mathcal{P}^b}\left(\frac{p_1}{p_1+p_4}, \frac{p_4}{p_1+p_4}\right)$ is arbitrarily close to ∞ , while, other than $c_{\mathcal{P}^b}\left(\frac{p_2}{p_2+p_3}, \frac{p_3}{p_2+p_3}\right)$, all other terms remain essentially constant. This then implies that $c_{\mathcal{P}^b}\left(\frac{p_2}{p_2+p_3}, \frac{p_3}{p_2+p_3}\right)$ drops by an arbitrarily large amount, which is not possible since it is positive by definition. Thus:

$$\liminf_{p \downarrow p_d} c_{\mathcal{P}^b}(p, 1-p) = L < c_{\mathcal{P}^b}(p_d, 1-p_d).$$

I next show that L is unbounded below, which causes the final contradiction since $c_{\mathcal{P}^b}(p, 1-p)$ cannot be negative. To do this, suppose that $L \geq 0$, define p_1, p_2, p_3 , and p_4 as above in this proof, and increase p_1 and decrease p_4 by an arbitrarily small $\delta > 0$, keeping p_2 and p_3 constant, so that $c_{\mathcal{P}^b}(p_1+p_2, p_3+p_4)$ is arbitrarily close to L . Then it is easy to see the contradiction using [Lemma 10](#) as in the previous paragraphs since $c_{\mathcal{P}^b}\left(\frac{p_1}{p_1+p_4}, \frac{p_4}{p_1+p_4}\right)$ is more than L , and all other terms remain essentially constant. ■

Proof of [Lemma 4](#).

Given a binary partition $\mathcal{P}^b = \{A_1, A_2\}$, define $c_{\mathcal{P}^b} : \cup_{j=1}^{\infty} \Delta^j \rightarrow \mathbb{R}$, where Δ^j is the j simplex, as in the statement of [Lemma 9](#), and suppose C satisfies [Axiom 1](#), [Axiom 2](#), and [Axiom 3](#). Remember [Lemma 9](#) tells me $c_{\mathcal{P}^b}(0, 1) = 0$, so I only need to show $c_{\mathcal{P}^b}(p, 1-p)$ is non-decreasing for small increases to $p \in (0, 1/2)$.

I proceed by assuming there is a $p_d \in (0, 1/2)$ such that $c_{\mathcal{P}^b}(p_d, 1-p_d)$ is decreasing for small increases in p_d and create a contradiction. Notice first that since [Lemma 3](#) shows $c_{\mathcal{P}^b}(p, 1-p)$ is continuous and [Lemma 9](#) shows $c_{\mathcal{P}^b}(0, 1) = 0$ that it must be that before any p where $c_{\mathcal{P}^b}(p, 1-p)$ is locally decreasing in p there must be a smaller p where $c_{\mathcal{P}^b}(p, 1-p)$ is locally increasing in p . Notice second that there must be infinitely many $p \in (0, 1/2)$ where $c_{\mathcal{P}^b}(p, 1-p)$ decreases for small increases to p because if $p_d \in (0, 1/2)$ is such that $c_{\mathcal{P}^b}(p_d, 1-p_d)$ decreases for small increases to p_d I can pick (p_1, p_2, p_2, p_4) such that:

$$p_1 + p_2 = p_d, \frac{p_1}{p_1+p_2} = p_d, \frac{p_3}{p_3+p_4} = p_d, \text{ so that } \frac{p_1}{p_1+p_4} < p_d,$$

and then notice [Lemma 10](#) tells us:

$$\begin{aligned} & c_{\mathcal{P}^b}(p_1, p_2, p_3, p_4) \\ &= c_{\mathcal{P}^b}(p_1+p_2, p_3+p_4) + (p_1+p_2)c_{\mathcal{P}^b}\left(\frac{p_1}{p_1+p_2}, \frac{p_2}{p_1+p_2}\right) + (p_3+p_4)c_{\mathcal{P}^b}\left(\frac{p_3}{p_3+p_4}, \frac{p_4}{p_3+p_4}\right) \end{aligned}$$

$$= c_{\mathcal{P}^b}(p_1 + p_4, p_2 + p_3) + (p_1 + p_4)c_{\mathcal{P}^b}\left(\frac{p_1}{p_1 + p_4}, \frac{p_4}{p_1 + p_4}\right) + (p_2 + p_3)c_{\mathcal{P}^b}\left(\frac{p_2}{p_2 + p_3}, \frac{p_3}{p_2 + p_3}\right),$$

and then consider increasing p_1 a small amount and decreasing p_4 by the same small amount, while keeping p_2 and p_3 constant, and notice this implies $c_{\mathcal{P}^b}(p, 1 - p)$ decreases for small increases to $p = p_1/(p_1 + p_4) < p_d$. This all means $c_{\mathcal{P}^b}(p, 1 - p)$ has dense local maxima and minima for p close to zero.

Next I show that the largest reduction in $c_{\mathcal{P}^b}(p, 1 - p)$ from an increase in p of any particular small $\epsilon > 0$ must be achieved at a $p > 1/4$. Pick $p_1 \leq 1/4$ such that $c_{\mathcal{P}^b}$ is decreasing there for small increases in p_1 . Given such a small $\epsilon > 0$, pick p_2 and p_3 so that $p_1 + p_2 + p_3 = 1$, and so:

$$\frac{p_3}{p_2 + p_3} = \frac{p_2 - \epsilon}{p_2 - \epsilon + p_3}.$$

Since ϵ is small and $p_1 \leq 1/4$, I know $p_1 < p_3 < p_2$. Pick $k \geq 0$ so:

$$k = c_{\mathcal{P}^b}\left(\frac{p_3}{p_2 + p_3}, 1 - \frac{p_3}{p_2 + p_3}\right) = c_{\mathcal{P}^b}\left(\frac{p_2 - \epsilon}{p_2 - \epsilon + p_3}, 1 - \frac{p_2 - \epsilon}{p_2 - \epsilon + p_3}\right).$$

[Lemma 9](#) and [Lemma 10](#) (or [Axiom 2](#)) tell us:

$$\begin{aligned} c_{\mathcal{P}^b}(p_1, p_2, p_3) &= c_{\mathcal{P}^b}(p_3, 1 - p_3) + (1 - p_3)c_{\mathcal{P}^b}\left(\frac{p_1}{p_1 + p_2}, \frac{p_2}{p_1 + p_2}\right) \\ &= c_{\mathcal{P}^b}(p_1, 1 - p_1) + (1 - p_1)c_{\mathcal{P}^b}\left(\frac{p_2}{p_2 + p_3}, \frac{p_3}{p_2 + p_3}\right). \end{aligned}$$

So, if I increase p_1 by ϵ and decrease p_2 by ϵ , the change in $c_{\mathcal{P}^b}(p_1, p_2, p_3)$ is:

$$\begin{aligned} &(1 - p_3) \left(c_{\mathcal{P}^b}\left(\frac{p_1 + \epsilon}{p_1 + p_2}, \frac{p_2 - \epsilon}{p_1 + p_2}\right) - c_{\mathcal{P}^b}\left(\frac{p_1}{p_1 + p_2}, \frac{p_2}{p_1 + p_2}\right) \right) \\ &= c_{\mathcal{P}^b}(p_1 + \epsilon, 1 - (p_1 + \epsilon)) - c_{\mathcal{P}^b}(p_1, 1 - p_1) - \epsilon k < 0. \end{aligned}$$

This implies:

$$\begin{aligned} &\frac{c_{\mathcal{P}^b}\left(\frac{p_1}{p_1 + p_2} + \frac{\epsilon}{p_1 + p_2}, \frac{p_2}{p_1 + p_2} - \frac{\epsilon}{p_1 + p_2}\right) - c_{\mathcal{P}^b}\left(\frac{p_1}{p_1 + p_2}, \frac{p_2}{p_1 + p_2}\right)}{\frac{\epsilon}{p_1 + p_2}} \\ &\leq \frac{c_{\mathcal{P}^b}(p_1 + \epsilon, 1 - (p_1 + \epsilon)) - c_{\mathcal{P}^b}(p_1, 1 - p_1)}{\epsilon} < 0 \end{aligned}$$

Thus, at

$$\frac{p_1}{p_1 + p_2} > p_1,$$

$c_{\mathcal{P}^b}$ is averaging a weakly steeper descent over a longer range, and thus there must be a point between

$$\frac{p_1}{p_1 + p_2} \text{ and } \frac{p_1 + \epsilon}{p_1 + p_2},$$

where the decrease of $c_{\mathcal{P}^b}$ over the next ϵ is as large as the decrease $c_{\mathcal{P}^b}(p_1 + \epsilon, 1 - (p_1 + \epsilon)) - c_{\mathcal{P}^b}(p_1, 1 - p_1)$. When p_1 is close to $1/4$, if I pick p_2 and p_3 as above, keeping our small ϵ in mind, I have:

$$\frac{p_1}{p_1 + p_2} > \frac{1}{4}.$$

$c_{\mathcal{P}^b}$ is a continuous function, so for all small $\epsilon > 0$, $f(p) = c_{\mathcal{P}^b}(p + \epsilon, 1 - (p + \epsilon)) - c_{\mathcal{P}^b}(p, 1 - p)$, defined for compact domain $p \in [0, \frac{1}{2} - \epsilon]$, is continuous, and has a minimizer (perhaps not unique) $p_s(\epsilon) \in (\frac{1}{4}, \frac{1}{2} - \epsilon)$, given what I just showed.

I am now ready to create the desired contradiction. I begin by finding a local maximum, denote it p_m , such that $p_m \in (0, 1/1000)$, and an $\epsilon \in (0, 1/1000)$, such that if $\delta \in [0, \epsilon]$, then:

$$c_{\mathcal{P}^b}(p_m, 1 - p_m) > c_{\mathcal{P}^b}(p_m + 4\delta, 1 - (p_m + 4\delta)).$$

Now, let $p_2 = p_s(\epsilon) + \epsilon > 1/4 + \epsilon$, and let:

$$p_3 = \frac{p_2}{1 - p_m} p_m < p_m, \text{ so that } \frac{p_3}{p_2 + p_3} = p_m.$$

Finally, let $p_1 = 1 - p_2 - p_3$, noticing $p_1 > 1/4$ so that:

$$\frac{p_3}{p_1 + p_3} + \frac{\epsilon}{p_1 + p_3 + \epsilon} < \frac{1}{2}.$$

[Lemma 10](#) tells us:

$$\begin{aligned} c_{\mathcal{P}^b}(p_1, p_2, p_3) &= c_{\mathcal{P}^b}(p_1, 1 - p_1) + (1 - p_1)c_{\mathcal{P}^b}\left(\frac{p_2}{p_2 + p_3}, \frac{p_3}{p_2 + p_3}\right) \\ &= c_{\mathcal{P}^b}(p_2, 1 - p_2) + (1 - p_2)c_{\mathcal{P}^b}\left(\frac{p_1}{p_1 + p_3}, \frac{p_3}{p_1 + p_3}\right). \end{aligned}$$

This means, since $p_2 + p_3 > 1/4$, if I increase p_3 by ϵ , and decrease p_2 by ϵ , holding p_1 constant,

and consider the change in $c_{\mathcal{P}^b}(p_1, p_2, p_3)$:

$$\begin{aligned}
0 &> (1 - p_1) \left(c_{\mathcal{P}^b} \left(\frac{p_3 + \epsilon}{p_2 + p_3}, \frac{p_2 - \epsilon}{p_2 + p_3} \right) - c_{\mathcal{P}^b} \left(\frac{p_3}{p_2 + p_3}, \frac{p_2}{p_2 + p_3} \right) \right) \\
&= c_{\mathcal{P}^b}(p_2 - \epsilon, 1 - (p_2 - \epsilon)) - c_{\mathcal{P}^b}(p_2, 1 - p_2) \\
&+ (p_1 + p_3 + \epsilon) c_{\mathcal{P}^b} \left(\frac{p_3 + \epsilon}{p_1 + p_3 + \epsilon}, \frac{p_1}{p_1 + p_3 + \epsilon} \right) - (p_1 + p_3) c_{\mathcal{P}^b} \left(\frac{p_3}{p_1 + p_3}, \frac{p_1}{p_1 + p_3} \right) \\
&\geq c_{\mathcal{P}^b}(p_2 - \epsilon, 1 - (p_2 - \epsilon)) - c_{\mathcal{P}^b}(p_2, 1 - p_2) \\
&+ (p_1 + p_3 + \epsilon) \left(c_{\mathcal{P}^b} \left(\frac{p_3 + \epsilon}{p_1 + p_3 + \epsilon}, \frac{p_1}{p_1 + p_3 + \epsilon} \right) - c_{\mathcal{P}^b} \left(\frac{p_3}{p_1 + p_3}, \frac{p_1}{p_1 + p_3} \right) \right) \\
&= c_{\mathcal{P}^b}(p_2 - \epsilon, 1 - (p_2 - \epsilon)) - c_{\mathcal{P}^b}(p_2, 1 - p_2) \\
&+ (p_1 + p_3 + \epsilon) \left(c_{\mathcal{P}^b} \left(\frac{p_3}{p_1 + p_3 + \epsilon} + \frac{\epsilon}{p_1 + p_3 + \epsilon}, \frac{p_1}{p_1 + p_3 + \epsilon} \right) - c_{\mathcal{P}^b} \left(\frac{p_3}{p_1 + p_3}, \frac{p_1}{p_1 + p_3} \right) \right).
\end{aligned}$$

This implies:

$$\begin{aligned}
0 &> \frac{c_{\mathcal{P}^b}(p_s(\epsilon) + \epsilon, 1 - (p_s(\epsilon) + \epsilon)) - c_{\mathcal{P}^b}(p_s(\epsilon), 1 - p_s(\epsilon))}{\epsilon} \\
&> \frac{c_{\mathcal{P}^b} \left(\frac{p_3}{p_1 + p_3 + \epsilon} + \frac{\epsilon}{p_1 + p_3 + \epsilon}, \frac{p_1}{p_1 + p_3 + \epsilon} \right) - c_{\mathcal{P}^b} \left(\frac{p_3}{p_1 + p_3}, \frac{p_1}{p_1 + p_3} \right)}{\frac{\epsilon}{p_1 + p_3 + \epsilon}}.
\end{aligned}$$

But remember, the way I picked $p_s(\epsilon)$ implies for all $\delta \in \left[\epsilon, \frac{\epsilon}{p_1 + p_3 + \epsilon} \right]$:

$$\begin{aligned}
&\frac{c_{\mathcal{P}^b}(p_s(\epsilon) + \epsilon, 1 - (p_s(\epsilon) + \epsilon)) - c_{\mathcal{P}^b}(p_s(\epsilon), 1 - p_s(\epsilon))}{\epsilon} \\
&\leq \frac{c_{\mathcal{P}^b} \left(\frac{p_3}{p_1 + p_3} + \delta, \frac{p_1}{p_1 + p_3} - \delta \right) - c_{\mathcal{P}^b} \left(\frac{p_3}{p_1 + p_3}, \frac{p_1}{p_1 + p_3} \right)}{\delta},
\end{aligned}$$

so letting $\delta = \frac{\epsilon}{p_1 + p_3 + \epsilon} \frac{p_1}{p_1 + p_3} \in \left[\epsilon, \frac{\epsilon}{p_1 + p_3 + \epsilon} \right]$:

$$\begin{aligned}
&\frac{c_{\mathcal{P}^b}(p_s(\epsilon) + \epsilon, 1 - (p_s(\epsilon) + \epsilon)) - c_{\mathcal{P}^b}(p_s(\epsilon), 1 - p_s(\epsilon))}{\epsilon} \\
&\leq \frac{c_{\mathcal{P}^b} \left(\frac{p_3}{p_1 + p_3} + \frac{\epsilon}{p_1 + p_3 + \epsilon} \frac{p_1}{p_1 + p_3}, \frac{p_1}{p_1 + p_3} - \frac{\epsilon}{p_1 + p_3 + \epsilon} \frac{p_1}{p_1 + p_3} \right) - c_{\mathcal{P}^b} \left(\frac{p_3}{p_1 + p_3}, \frac{p_1}{p_1 + p_3} \right)}{\frac{\epsilon}{p_1 + p_3 + \epsilon} \frac{p_1}{p_1 + p_3}}
\end{aligned}$$

$$\begin{aligned}
&= \frac{c_{\mathcal{P}^b} \left(\frac{p_3}{p_1 + p_3 + \epsilon} + \frac{\epsilon}{p_1 + p_3 + \epsilon}, \frac{p_1 + \epsilon}{p_1 + p_3 + \epsilon} - \frac{\epsilon}{p_1 + p_3 + \epsilon} \right) - c_{\mathcal{P}^b} \left(\frac{p_3}{p_1 + p_3}, \frac{p_1}{p_1 + p_3} \right)}{\frac{\epsilon}{p_1 + p_3 + \epsilon} \frac{p_1}{p_1 + p_3}} \\
&< \frac{c_{\mathcal{P}^b} \left(\frac{p_3}{p_1 + p_3 + \epsilon} + \frac{\epsilon}{p_1 + p_3 + \epsilon}, \frac{p_1}{p_1 + p_3 + \epsilon} \right) - c_{\mathcal{P}^b} \left(\frac{p_3}{p_1 + p_3}, \frac{p_1}{p_1 + p_3} \right)}{\frac{\epsilon}{p_1 + p_3 + \epsilon}},
\end{aligned}$$

which established the desired contradiction. ■

Proof of Lemma 5. Assume C satisfies [Axiom 1](#), [Axiom 2](#), and [Axiom 3](#). Given learning strategy invariant partition $\mathcal{P} = \{A_1, \dots, A_m\}$ pick any binary partition \mathcal{P}^b coarser than \mathcal{P} and define $c_{\mathcal{P}^b} : \cup_{j=1}^{\infty} \Delta^j \rightarrow \mathbb{R}$, where Δ^j is the j simplex, as in the statement of [Lemma 9](#) so that, by [Lemma 1](#), $C(\mathcal{P}, \mu) = c_{\mathcal{P}^b}(\mu(A_1), \dots, \mu(A_m))$.

I begin by showing that if there is a $p \in (0, \frac{1}{2}]$ such that $c_{\mathcal{P}^b}(p, 1-p) = 0$, then $c_{\mathcal{P}^b}(p, 1-p) = 0 \forall p \in (0, \frac{1}{2}]$. Fix $p \in (0, \frac{1}{2}]$, to be the largest number less than $\frac{1}{2}$ such that $c_{\mathcal{P}^b}(p, 1-p) = 0$, let $p_1 = p_2 = p$, and let $p_3 = 1 - p_1 - p_2$. [Lemma 9](#) and [Lemma 10](#) tell us that: $c_{\mathcal{P}^b}(p_1, p_2, p_3) =$

$$c_{\mathcal{P}^b}(p_1, 1-p) + (1-p_1)c_{\mathcal{P}^b} \left(\frac{p_2}{p_2 + p_3}, \frac{p_3}{p_2 + p_3} \right) = c_{\mathcal{P}^b}(p_3, 1-p_3) + (1-p_3)c_{\mathcal{P}^b} \left(\frac{p_1}{p_1 + p_2}, \frac{p_2}{p_1 + p_2} \right).$$

Now suppose $p_1 > 0$ and that I decrease p_1 and increase p_2 by the same arbitrarily small $\epsilon > 0$. The result in [Lemma 4](#) creates a contradiction, however, since $\frac{p_2}{p_2 + p_3} > p_1$, so:

$$c_{\mathcal{P}^b}(p_1 - \epsilon, 1 - (p_1 - \epsilon)) + (1 - (p_1 - \epsilon))c_{\mathcal{P}^b} \left(\frac{p_2 + \epsilon}{p_2 + \epsilon + p_3}, \frac{p_3}{p_2 + \epsilon + p_3} \right)$$

would be increasing in ϵ , while

$$c_{\mathcal{P}^b}(p_3, 1 - p_3) + (1 - p_3)c_{\mathcal{P}^b} \left(\frac{p_1 - \epsilon}{p_1 + p_2}, \frac{p_2 + \epsilon}{p_1 + p_2} \right)$$

would be decreasing in ϵ . Thus p_1 cannot be strictly positive, and it must be that $p = 0$.

This all means that if $\exists p \in (0, \frac{1}{2}]$ such that $c_{\mathcal{P}^b}(p, 1-p) = 0$, then:

$$C(\mathcal{P}, \mu) = c_{\mathcal{P}^b}(\mu(A_1), \dots, \mu(A_m)) = 0 = 0\mathcal{H}(\mathcal{P}, \mu).$$

For the rest of the proof I assume $c_{\mathcal{P}^b}(p, 1-p) > 0 \forall p \in (0, \frac{1}{2}]$. Define h so that for $n \in \mathbb{N}$, $h(n) \equiv c_{\mathcal{P}^b}(1/n, \dots, 1/n, 0)$. Since I assumed, $c_{\mathcal{P}^b}(p, 1-p) > 0 \forall p \in (0, \frac{1}{2}]$, $h(2) > h(1) = 0$, and in general $h(n) > 0$ if $n > 1$. It is also easy to show $h(n+1) > h(n)$ for all $n \geq 2$ using [Lemma 10](#)

and [Lemma 4](#):

$$\begin{aligned}
h(n) &= c_{\mathcal{P}^b}(1/n, \dots, 1/n, 0) \\
&= c_{\mathcal{P}^b}(1/n, \dots, 1/n) + \left(\frac{1}{n}\right) c_{\mathcal{P}^b}\left(\frac{1/n}{1/n}, \frac{0}{1/n}\right) \\
&< c_{\mathcal{P}^b}(1/n, \dots, 1/n) + \left(\frac{1}{n}\right) c_{\mathcal{P}^b}\left(\frac{1}{\frac{1}{n}}, \frac{1}{\frac{1}{n}}\right) \\
&= c_{\mathcal{P}^b}(1/n, \dots, 1/n, 1/n, 1/(n+1), 1/(n(n+1))) = c_{\mathcal{P}^b}(1/n, \dots, 1/n, 1/(n+1), 1/n, 1/(n(n+1))) \\
&= c_{\mathcal{P}^b}(1/n, \dots, 1/n, 1/(n+1), (1/n) + 1/(n(n+1))) + \frac{n+2}{n(n+1)} c_{\mathcal{P}^b}\left(\frac{\frac{1}{n}}{\frac{n+2}{n(n+1)}}, \frac{\frac{1}{n(n+1)}}{\frac{n+2}{n(n+1)}}\right) \\
&\leq c_{\mathcal{P}^b}(1/n, \dots, 1/n, 1/(n+1), (1/n) + 1/(n(n+1))) + \frac{n+2}{n(n+1)} c_{\mathcal{P}^b}\left(\frac{\frac{1}{n+1}}{\frac{n+2}{n(n+1)}}, \frac{\frac{2}{n(n+1)}}{\frac{n+2}{n(n+1)}}\right) \\
&\leq \dots \leq c_{\mathcal{P}^b}(1/(n+1), \dots, 1/(n+1), 0) = h(n+1)
\end{aligned}$$

The rest of the proof follows the work of [Shannon \(1948\)](#) closely. Notice $h(s^r) = r \cdot h(s)$, which is reminiscent of logarithms, and is some nice foreshadowing for the rest of the proof. Given arbitrarily small $\epsilon > 0$, and integers $s > 1$ and $t > 1$, pick n and r so that $2/n < \epsilon$, and $s^r \leq t^n < s^{r+1}$. So:

$$r \log(s) \leq n \log(t) < (r+1) \log(s) \implies \frac{r}{n} \leq \frac{\log(t)}{\log(s)} < \frac{r+1}{n} \implies \left| \frac{r}{n} - \frac{\log(t)}{\log(s)} \right| < \frac{1}{n}.$$

The work I did above then tells us:

$$\begin{aligned}
h(s^r) \leq h(t^n) \leq h(s^{r+1}) &\implies r \cdot h(s) \leq n \cdot h(t) \leq (r+1)h(s) \\
\implies \frac{r}{n} \leq \frac{h(t)}{h(s)} \leq \frac{r+1}{n} &\implies \left| \frac{r}{n} - \frac{h(t)}{h(s)} \right| \leq \frac{1}{n}.
\end{aligned}$$

All of this tells us:

$$\left| \frac{h(t)}{h(s)} - \frac{\log(t)}{\log(s)} \right| < \epsilon,$$

which can be shown to be true $\forall \epsilon > 0$, and thus $h(n) = \lambda \log(n)$, where λ must be a positive constant.

Let $p_k = \mu(A_k)$ for each $A_k \in \mathcal{P}$. Suppose, for now, that each p_k is a rational number. Then

there exists integers n_1, \dots, n_m , such that for all $k \in \{1, \dots, m\}$ I have:

$$p_k = \frac{n_k}{\sum_{j=1}^m n_j}.$$

The interpretation is that I have a uniform distribution over $\sum_j n_j$ equally likely states, and the chance of the event which happens with probability p_k is the chance of one of the n_k associated states occurring. Then using the definition of learning strategy invariance:

$$\begin{aligned} c_{\mathcal{P}^b} \left(\frac{1}{\sum_j n_j}, \dots, \frac{1}{\sum_j n_j} \right) &= h \left(\sum_{j=1}^m n_j \right) = \lambda \log \left(\sum_{j=1}^m n_j \right) = c_{\mathcal{P}^b}(p_1, \dots, p_m) + \sum_{j=1}^m p_j \lambda \log(n_j), \\ \implies c_{\mathcal{P}^b}(p_1, \dots, p_m) &= \lambda \log \left(\sum_{j=1}^m n_j \right) - \sum_{j=1}^m p_j \lambda \log(n_j) \\ &= \sum_{k=1}^m \left(p_k \lambda \log \left(\sum_{j=1}^m n_j \right) \right) - \sum_{j=1}^m p_j \lambda \log(n_j) \\ &= - \sum_{k=1}^m p_k \lambda \log \left(\frac{n_k}{\sum_j n_j} \right) = -\lambda \sum_{k=1}^m p_k \log(p_k) = \lambda \mathcal{H}(\mathcal{P}, \mu), \end{aligned}$$

where \mathcal{H} is defined as in equation (14). If any of the p_i are irrational, then the density of the rationals and [Lemma 3](#) can be used to get the same result. Thus:

$$C(\mathcal{P}, \mu) = c_{\mathcal{P}^b}(\mu(A_1), \dots, \mu(A_m)) = \lambda \mathcal{H}(\mathcal{P}, \mu). \blacksquare$$

Mutual Information

Consider two partitions \mathcal{P}_1 and \mathcal{P}_2 . Given some probability measure μ , define the **mutual information** between \mathcal{P}_1 and \mathcal{P}_2 , denoted $I(\mathcal{P}_1, \mathcal{P}_2, \mu)$, to be:

$$I(\mathcal{P}_1, \mathcal{P}_2, \mu) = \sum_{a_1 \in \mathcal{P}_1} \sum_{a_2 \in \mathcal{P}_2} \mu(a_1 \cap a_2) \log \left(\frac{\mu(a_1 \cap a_2)}{\mu(a_1)\mu(a_2)} \right)$$

Then, as is well known in the literature:

$$\mathcal{H}(\times\{\mathcal{P}_i\}_{i=1}^2, \mu) = \mathcal{H}(\mathcal{P}_1, \mu) + \mathcal{H}(\mathcal{P}_2, \mu) - I(\mathcal{P}_1, \mathcal{P}_2, \mu)$$

$$\begin{aligned}
&= \mathbb{E}[\mathcal{H}(\mathcal{P}_1, \mu(\cdot|\mathcal{P}_2(\omega)))] + I(\mathcal{P}_1, \mathcal{P}_2, \mu) + \mathbb{E}[\mathcal{H}(\mathcal{P}_2, \mu(\cdot|\mathcal{P}_1(\omega)))] \\
&\quad \parallel \qquad \qquad \qquad \parallel \\
&\quad \mathcal{H}(\mathcal{P}_1, \mu) - I(\mathcal{P}_1, \mathcal{P}_2, \mu) \qquad \qquad \mathcal{H}(\mathcal{P}_2, \mu) - I(\mathcal{P}_1, \mathcal{P}_2, \mu) \\
&= \mathcal{H}(\mathcal{P}_1, \mu) + \mathbb{E}[\mathcal{H}(\mathcal{P}_2, \mu(\cdot|\mathcal{P}_1(\omega)))] = \mathcal{H}(\mathcal{P}_2, \mu) + \mathbb{E}[\mathcal{H}(\mathcal{P}_1, \mu(\cdot|\mathcal{P}_2(\omega)))]
\end{aligned}$$

and note that the strict concavity of \mathcal{H} means that $I(\mathcal{P}_1, \mathcal{P}_2, \mu) \geq 0$.

Mutual information can be thought of as the information that is double counted if one were to compute the total uncertainty about the outcome of \mathcal{P}_1 and \mathcal{P}_2 by simply adding up the uncertainty about the outcome of \mathcal{P}_1 and the uncertainty about the outcome of \mathcal{P}_2 . When the mutual information increases and the individual uncertainty about the outcome of \mathcal{P}_1 and the outcome of \mathcal{P}_2 are held constant the total uncertainty about the outcome of \mathcal{P}_1 and \mathcal{P}_2 decreases because the amount that remains to be learned after observing one of the outcomes of either \mathcal{P}_1 or \mathcal{P}_2 decreases.

Mutual information can be acquired by learning the value of either \mathcal{P}_1 or \mathcal{P}_2 . When I think of an agent that is trying to acquire information in an efficient fashion, I should always envision them acquiring mutual information from the cheapest attribute, by learning about whichever of \mathcal{P}_1 and \mathcal{P}_2 has the lowest associated multiplier. This logic is formalized by the result in [Lemma 11](#).

Lemma 11. If C satisfies our four axioms, and $S^b = \{\mathcal{P}_1^b, \dots, \mathcal{P}_i^b, \mathcal{P}_{i+1}^b, \dots, \mathcal{P}_m^b\}$ and $\tilde{S}^b = \{\mathcal{P}_1^b, \dots, \mathcal{P}_{i+1}^b, \mathcal{P}_i^b, \dots, \mathcal{P}_m^b\}$ are two binary learning strategies such that \mathcal{P}_i^b and \mathcal{P}_{i+1}^b 's associated multipliers are ordered $\lambda_i \geq \lambda_{i+1}$, then for all probability measures μ :

$$C(S^b, \mu) \geq C(\tilde{S}^b, \mu).$$

Proof. For all realizations of $\cap_{j=1}^{i-1} \mathcal{P}_j^b(\omega)$:

$$\begin{aligned}
C((\mathcal{P}_i^b, \mathcal{P}_{i+1}^b), \mu(\cdot|\cap_{j=1}^{i-1} \mathcal{P}_j^b(\omega))) &= \lambda_i \mathcal{H}(\mathcal{P}_i^b, \mu(\cdot|\cap_{j=1}^{i-1} \mathcal{P}_j^b(\omega))) + \lambda_{i+1} \mathbb{E}[\mathcal{H}(\mathcal{P}_{i+1}^b, \mu(\cdot|\cap_{j=1}^i \mathcal{P}_j^b(\omega)))] \\
&= \lambda_i \mathcal{H}(\mathcal{P}_i^b, \mu(\cdot|\cap_{j=1}^{i-1} \mathcal{P}_j^b(\omega))) + \lambda_{i+1} \left(\mathcal{H}(\mathcal{P}_{i+1}^b, \mu(\cdot|\cap_{j=1}^{i-1} \mathcal{P}_j^b(\omega))) - I(\mathcal{P}_i^b, \mathcal{P}_{i+1}^b, \mu(\cdot|\cap_{j=1}^{i-1} \mathcal{P}_j^b(\omega))) \right) \\
&\geq \lambda_i \left(\mathcal{H}(\mathcal{P}_i^b, \mu(\cdot|\cap_{j=1}^{i-1} \mathcal{P}_j^b(\omega))) - I(\mathcal{P}_i^b, \mathcal{P}_{i+1}^b, \mu(\cdot|\cap_{j=1}^{i-1} \mathcal{P}_j^b(\omega))) \right) + \lambda_{i+1} \mathcal{H}(\mathcal{P}_{i+1}^b, \mu(\cdot|\cap_{j=1}^{i-1} \mathcal{P}_j^b(\omega))) \\
&= \lambda_{i+1} \mathcal{H}(\mathcal{P}_{i+1}^b, \mu(\cdot|\cap_{j=1}^{i-1} \mathcal{P}_j^b(\omega))) + \lambda_i \mathbb{E}[\mathcal{H}(\mathcal{P}_i^b, \mu(\cdot|(\cap_{j=1}^{i-1} \mathcal{P}_j^b(\omega)) \cap \mathcal{P}_{i+1}^b(\omega)))] \\
&= C((\mathcal{P}_{i+1}^b, \mathcal{P}_i^b), \mu(\cdot|\cap_{j=1}^{i-1} \mathcal{P}_j^b(\omega))).
\end{aligned}$$

It is thus always weakly cheaper in expectation to have \mathcal{P}_{i+1} before \mathcal{P}_i since switching their

order does not change the expected cost of implementing the binary partitions before or after the pair. ■

Proof of Theorem 1. Given some probability measure μ , suppose S^b is a binary learning strategy such that $\sigma(S^b) = \mathcal{F}$, and

$$C(S^b, \mu) = \min_{S^b \in \mathcal{S}^b(\Omega)} C(S^b, \mu).$$

I may assume that if \mathcal{P}_i^b and \mathcal{P}_{i+1}^b are in S^b with associated multipliers λ_i and λ_{i+1} , that $\lambda_i \leq \lambda_{i+1}$. If not, then their order can be reversed and the resultant strategy is weakly less costly, as is shown in [Lemma 11](#).

If for any $j \in \{1, \dots, M\}$, multiplier λ_j 's associated binary partitions $\mathcal{P}_i^b, \dots, \mathcal{P}_{i+k}^b$ in S^b are such that $\sigma(\mathcal{P}_i^b, \dots, \mathcal{P}_{i+k}^b) \neq \sigma(\mathcal{P}_{\lambda_j}^b)$, then there are binary partitions $\mathcal{P}_{m+1}^b, \dots, \mathcal{P}_{m+l}^b$ with associated multiplier λ_j , such that $\sigma(\mathcal{P}_i^b, \dots, \mathcal{P}_{i+k}^b, \mathcal{P}_{m+1}^b, \dots, \mathcal{P}_{m+l}^b) = \sigma(\mathcal{P}_{\lambda_j}^b)$. $\mathcal{P}_{m+1}^b, \dots, \mathcal{P}_{m+l}^b$ can be appended to the end of S^b , and the resultant strategy \tilde{S}^b is also such that:

$$C(\tilde{S}^b, \mu) = \min_{S^b \in \mathcal{S}^b(\Omega)} C(S, \mu).$$

This is true since each appended binary partition has an expected cost of zero, since $\sigma(S^b) = \mathcal{F}$. [Lemma 11](#) then implies that if I reorder \tilde{S}^b so that the new learning strategy \hat{S}^b 's binary partitions are ordered by their multipliers, then:

$$C(\hat{S}^b, \mu) = \min_{S^b \in \mathcal{S}^b(\Omega)} C(S, \mu).$$

I can thus assume that S^b is such that for any $j \in \{1, \dots, M\}$ multiplier λ_j 's associated binary partitions $\mathcal{P}_i^b, \dots, \mathcal{P}_{i+k}^b$ in S^b are such that $\sigma(\mathcal{P}_i^b, \dots, \mathcal{P}_{i+k}^b) = \sigma(\mathcal{P}_{\lambda_j}^b)$.

For each $j \in \{1, \dots, M\}$ I thus have that if all binary partitions $\mathcal{P}_i^b, \dots, \mathcal{P}_{i+k}^b$ in S^b with multiplier λ_j are taken together that:

$$\begin{aligned} \mathbb{E}[C((\mathcal{P}_i^b, \dots, \mathcal{P}_{i+k}^b), \mu(\cdot | \cap_{t=1}^{i-1} \mathcal{P}_t^b(\omega)))] &= \mathbb{E}\left[\sum_{l=i}^{i+k} \lambda_j \mathcal{H}(\mathcal{P}_l^b, \mu(\cdot | \cap_{t=1}^{l-1} \mathcal{P}_t^b(\omega)))\right] \\ &= \mathbb{E}[\lambda_j \mathcal{H}(\mathcal{P}_{\lambda_j}, \mu(\cdot | \cap_{t=1}^{i-1} \mathcal{P}_t^b(\omega)))] = \mathbb{E}[\lambda_j \mathcal{H}(\mathcal{P}_{\lambda_j}, \mu(\cdot | \cap_{t=1}^{j-1} \mathcal{P}_{\lambda_t}(\omega)))]. \end{aligned}$$

Where the second equality holds due to the properties of \mathcal{H} . This procedure can be carried out for

all μ . Thus:

$$\begin{aligned} C(S^b, \mu) &= \min_{S^b \in \mathcal{S}^b(\Omega)} C(S, \mu). \\ &= \lambda_1 \mathcal{H}(\mathcal{P}_{\lambda_1}, \mu) + \mathbb{E} \left[\lambda_2 \mathcal{H}(\mathcal{P}_{\lambda_2}, \mu(\cdot | \mathcal{P}_{\lambda_1}(\omega))) + \dots + \lambda_M \mathcal{H}(\mathcal{P}_{\lambda_M}, \mu(\cdot | \cap_{i=1}^{M-1} \mathcal{P}_{\lambda_i}(\omega))) \right]. \end{aligned}$$

This is equivalent to the equation in the statement of the theorem due to the definition of the attributes. ■

Appendix 2

Proof of Lemma 6. In Lemma 6, I show that I can rewrite the agent's problem in terms of selecting the choice probabilities described in equations (4), (5), and (6). To do this, I first establish several other lemmas. In Lemma 12, I show that: $\min_{S \in \mathcal{S}^b(\Omega)} C(S, \mu)$ is a strictly concave function of μ . This is a commonly known property of Shannon Entropy, but needs to be established for MASE. This implies that \mathbf{C} (defined in equation (1)) is strictly convex. I then show, in Lemma 13, that, given the convexity of \mathbf{C} , any selected action is associated with a particular posterior probability. This is desirable because it allows us to reduce the strategies considered to recommendation strategies. That is, I am able to focus on signals that are simply a recommendation of an option. In Lemma 14, I show that I may rewrite the cost function in terms of the choice probabilities in equations (4), (5), and (6).

Lemma 12. If C satisfies all four axioms then $\min_{S \in \mathcal{S}^b(\Omega)} C(S, \mu)$ is a strictly concave function of μ . Namely, if there are probability measures μ_a , and μ_b , such that, $\forall \omega \in \Omega$, $\mu(\omega) = \alpha \mu_a(\omega) + (1 - \alpha) \mu_b(\omega)$ for some $\alpha \in (0, 1)$, and $\mu_a \neq \mu_b$, then:

$$\min_{S \in \mathcal{S}^b(\Omega)} C(S, \mu) > \alpha \left(\min_{S \in \mathcal{S}^b(\Omega)} C(S, \mu_a) \right) + (1 - \alpha) \left(\min_{S \in \mathcal{S}^b(\Omega)} C(S, \mu_b) \right).$$

Proof. For each such μ_a , μ_b , $\alpha \in (0, 1)$, and μ , the strict concavity of Shannon Entropy (Matějka & McKay, 2015; Caplin et al., 2017) implies:

$$\mathcal{H}(\mathcal{A}_1, \mu) \geq \alpha \mathcal{H}(\mathcal{A}_1, \mu_a) + (1 - \alpha) \mathcal{H}(\mathcal{A}_1, \mu_b).$$

Define a random variable X that takes value 1 with chance α , and takes value 0 with chance $1 - \alpha$, so that a draw from μ is equivalent to a draw of X , and then a draw according to the probability measure $X\mu_a + (1 - X)\mu_b$. For each $i \in \{2, \dots, M\}$ and probability measure ν :

$\mathcal{P}_{\lambda_i} \times \{0, 1\} \rightarrow [0, 1]$, define:

$$\mathcal{H}(X, \nu) = \sum_X \nu(x) \log(\nu(x)), \quad \mathcal{H}(\mathcal{P}_{\lambda_i}, X, \nu) = \sum_{A \in \mathcal{P}_{\lambda_i}} \sum_X \nu(A, x) \log(\nu(A, x)).$$

Then, for each such $\mu_a, \mu_b, \alpha \in (0, 1)$, and μ such that $\mu = \alpha\mu_a + (1 - \alpha)\mu_b$, and $i \in \{2, \dots, M\}$, the properties of Shannon Entropy tell us:

$$\begin{aligned} \mathbb{E}\left[\mathcal{H}(\mathcal{P}_{\lambda_i}, X, \mu(\cdot | \cap_{j=1}^{i-1} \mathcal{P}_{\lambda_j}(\omega)))\right] &= \mathbb{E}\left[\mathcal{H}(\mathcal{P}_{\lambda_i}, \mu(\cdot | \cap_{j=1}^{i-1} \mathcal{P}_{\lambda_j}(\omega)))\right] + \mathbb{E}\left[\mathcal{H}(X, \mu(\cdot | \cap_{j=1}^i \mathcal{P}_{\lambda_j}(\omega)))\right], \\ \mathbb{E}\left[\mathcal{H}(\mathcal{P}_{\lambda_i}, X, \mu(\cdot | \cap_{j=1}^{i-1} \mathcal{P}_{\lambda_j}(\omega)))\right] &= \mathbb{E}\left[\mathcal{H}(X, \mu(\cdot | \cap_{j=1}^{i-1} \mathcal{P}_{\lambda_j}(\omega)))\right] + \mathbb{E}\left[\mathcal{H}(\mathcal{P}_{\lambda_i}, \mu(\cdot | \cap_{j=1}^{i-1} \mathcal{P}_{\lambda_j}(\omega), X))\right], \\ \implies \mathbb{E}\left[\mathcal{H}(\mathcal{P}_{\lambda_i}, \mu(\cdot | \cap_{j=1}^{i-1} \mathcal{P}_{\lambda_j}(\omega)))\right] &= \mathbb{E}\left[\mathcal{H}(\mathcal{P}_{\lambda_i}, \mu(\cdot | \cap_{j=1}^{i-1} \mathcal{P}_{\lambda_j}(\omega), X))\right] \\ &\quad + \mathbb{E}\left[\mathcal{H}(X, \mu(\cdot | \cap_{j=1}^{i-1} \mathcal{P}_{\lambda_j}(\omega)))\right] - \mathbb{E}\left[\mathcal{H}(X, \mu(\cdot | \cap_{j=1}^i \mathcal{P}_{\lambda_j}(\omega)))\right] \\ &\geq \mathbb{E}\left[\mathcal{H}(\mathcal{P}_{\lambda_i}, \mu(\cdot | \cap_{j=1}^{i-1} \mathcal{P}_{\lambda_j}(\omega), X))\right] \\ &= \mathbb{E}\left[\alpha \mathcal{H}(\mathcal{P}_{\lambda_i}, \mu_a(\cdot | \cap_{j=1}^{i-1} \mathcal{P}_{\lambda_j}(\omega))) + (1 - \alpha) \mathcal{H}(\mathcal{P}_{\lambda_i}, \mu_b(\cdot | \cap_{j=1}^{i-1} \mathcal{P}_{\lambda_j}(\omega)))\right]. \end{aligned}$$

The above inequality is strict for at least one $i \in \{2, \dots, M\}$ if the inequality from the previous paragraph is not strict, since $\mu_a \neq \mu_b$ and \mathcal{H} is strictly concave. The desired result thus follows from [Theorem 1](#) and the definition of the attributes. ■

Lemma 13. If action $n \in \mathcal{N}$ is selected with positive probability, $\Pr(n) > 0$, as the outcome of information strategy F which is a solution to (2) subject to (3), then there exists a posterior belief B_n such that $F(\omega|s) = B_n$ with probability one whenever n is selected.

Proof. It is impossible that there are two distinct sets of signals \mathcal{S}_n^1 and \mathcal{S}_n^2 which are observed with strictly positive probability, both of which lead to the selection of n , and induce different posteriors: $F(\omega|s_1) \neq F(\omega|s_2)$ for all $s_1 \in \mathcal{S}_n^1$ and $s_2 \in \mathcal{S}_n^2$. $\min_{S \in \mathcal{S}^b(\Omega)} C(S, \mu)$ is strictly concave in μ , as shown in [Lemma 12](#), so the agent could thus do better by replacing their original information strategy F with a new information strategy \tilde{F} which is identical to F except the signals in \mathcal{S}_n^1 and \mathcal{S}_n^2 are replaced by s_0 : $\forall \omega \in \Omega$ let $\tilde{F}(s_0|\omega) = \int_{s \in \mathcal{S}_n^1} F(s|\omega) + \int_{s \in \mathcal{S}_n^2} F(s|\omega)$. This is true because payoffs are linear, and the law of iterated expectations implies the agent still picks n after s_0 is realized

since $\forall \nu \in \mathcal{N}$:

$$\begin{aligned}
\mathbb{E}_{\tilde{F}}[\mathbf{v}_n(\omega)|s_0] &= \frac{\sum_{\omega \in \Omega} \int_{s \in \mathcal{S}_n^1} F(s|\omega) \mu(\omega)}{\sum_{\omega \in \Omega} \left(\int_{s \in \mathcal{S}_n^1} F(s|\omega) \mu(\omega) + \int_{s \in \mathcal{S}_n^2} F(s|\omega) \mu(\omega) \right)} \mathbb{E}_F[\mathbf{v}_n(\omega)|s \in \mathcal{S}_n^1] \\
&+ \frac{\sum_{\omega \in \Omega} \int_{s \in \mathcal{S}_n^2} F(s|\omega) \mu(\omega)}{\sum_{\omega \in \Omega} \left(\int_{s \in \mathcal{S}_n^1} F(s|\omega) \mu(\omega) + \int_{s \in \mathcal{S}_n^2} F(s|\omega) \mu(\omega) \right)} \mathbb{E}_F[\mathbf{v}_n(\omega)|s \in \mathcal{S}_n^2] \\
&\geq \frac{\sum_{\omega \in \Omega} \int_{s \in \mathcal{S}_n^1} F(s|\omega) \mu(\omega)}{\sum_{\omega \in \Omega} \left(\int_{s \in \mathcal{S}_n^1} F(s|\omega) \mu(\omega) + \int_{s \in \mathcal{S}_n^2} F(s|\omega) \mu(\omega) \right)} \mathbb{E}_F[\mathbf{v}_\nu(\omega)|s \in \mathcal{S}_n^1] \\
&+ \frac{\sum_{\omega \in \Omega} \int_{s \in \mathcal{S}_n^2} F(s|\omega) \mu(\omega)}{\sum_{\omega \in \Omega} \left(\int_{s \in \mathcal{S}_n^1} F(s|\omega) \mu(\omega) + \int_{s \in \mathcal{S}_n^2} F(s|\omega) \mu(\omega) \right)} \mathbb{E}_F[\mathbf{v}_\nu(\omega)|s \in \mathcal{S}_n^2] = \mathbb{E}_{\tilde{F}}[\mathbf{v}_\nu(\omega)|s_0]. \blacksquare
\end{aligned}$$

Lemma 14. The cost of information for a given strategy in equation (2) can be written:

$$\begin{aligned}
\mathbf{C}(F(s, \omega), \mu) &= \mathbf{C}(\mathbb{P}, \mu) \\
&\equiv \sum_{\omega \in \Omega} \mu(\omega) \sum_{n \in \mathcal{N}} \left(-\lambda_1 \Pr(n) \log(\Pr(n)) - (\lambda_2 - \lambda_1) \Pr(n|\mathcal{A}_1(\omega)) \log(\Pr(n|\mathcal{A}_1(\omega))) \right. \\
&\quad \left. - (\lambda_3 - \lambda_2) \Pr(n|\mathcal{A}_1(\omega) \cap \mathcal{A}_2(\omega)) \log(\Pr(n|\mathcal{A}_1(\omega) \cap \mathcal{A}_2(\omega))) \right. \\
&\quad \left. - \dots - (\lambda_M - \lambda_{M-1}) \Pr(n|\cap_{i=1}^{M-1} \mathcal{A}_i(\omega)) \log(\Pr(n|\cap_{i=1}^{M-1} \mathcal{A}_i(\omega))) + \lambda_M \Pr(n|\omega) \log(\Pr(n|\omega)) \right).
\end{aligned}$$

Proof. Let $\mathcal{P}_s = (\mathcal{S}_1, \dots, \mathcal{S}_N)$ denote a partition of the space of signals the agent may receive such that for each option n with $\Pr(n) > 0$ each signal that results in the agent selecting option n is in \mathcal{S}_n , and then as shown in Lemma 13, with probability one the s drawn from \mathcal{S}_n results in a particular posterior. I then have (using the properties of \mathcal{H}):

$$\mathbf{C}(F(s, \omega), \mu) \equiv \mathbb{E} \left[\min_{S \in \mathcal{S}^b(\Omega)} C(S, \mu) - \min_{S \in \mathcal{S}^b(\Omega)} C(S, \mu(\cdot|s)) \right]$$

$$= \mathbb{E} \left[\lambda_1 \left(\mathcal{H}(\mathcal{A}_1, \mu) - \mathcal{H}(\mathcal{A}_1, \mu(\cdot|s)) \right) \right] \quad (11)$$

$$+ \dots + \lambda_M \left(\mathcal{H}(\mathcal{A}_M, \mu(\cdot| \cap_{i=1}^{M-1} \mathcal{A}_i(\omega))) - \mathcal{H}(\mathcal{A}_M, \mu(\cdot| \cap_{i=1}^{M-1} \mathcal{A}_i(\omega), s)) \right) \Big]$$

$$= \mathbb{E} \left[\lambda_1 \left(\mathcal{H}(\mathcal{P}_s, F(s)) - \mathcal{H}(\mathcal{P}_s, F(s|\mathcal{A}_1(\omega))) \right) \right] \quad (12)$$

$$+ \dots + \lambda_M \left(\mathcal{H}(\mathcal{P}_s, F(s| \cap_{i=1}^{M-1} \mathcal{A}_i(\omega))) - \mathcal{H}(\mathcal{P}_s, F(s| \cap_{i=1}^M \mathcal{A}_i(\omega))) \right) \Big]$$

$$= \mathbb{E} \left[\lambda_1 \mathcal{H}(\mathcal{P}_s, F(s)) + (\lambda_2 - \lambda_1) \mathcal{H}(\mathcal{P}_s, F(s|\mathcal{A}_1(\omega))) \right]$$

$$+ \dots + (\lambda_M - \lambda_{M-1}) \mathcal{H}(\mathcal{P}_s, F(s| \cap_{i=1}^{M-1} \mathcal{A}_i(\omega))) - \lambda_M \mathcal{H}(\mathcal{P}_s, F(s| \cap_{i=1}^M \mathcal{P}_{\lambda_i}(\omega))) \Big]$$

$$= \sum_{\omega \in \Omega} \mu(\omega) \sum_{n \in \mathcal{N}} \left(-\lambda_1 \Pr(n) \log(\Pr(n)) - (\lambda_2 - \lambda_1) \Pr(n|\mathcal{A}_1(\omega)) \log(\Pr(n|\mathcal{A}_1(\omega))) \right)$$

$$- (\lambda_3 - \lambda_2) \Pr(n|\mathcal{A}_1(\omega) \cap \mathcal{A}_2(\omega)) \log(\Pr(n|\mathcal{A}_1(\omega) \cap \mathcal{A}_2(\omega)))$$

$$- \dots - (\lambda_M - \lambda_{M-1}) \Pr(n| \cap_{i=1}^{M-1} \mathcal{A}_i(\omega)) \log(\Pr(n| \cap_{i=1}^{M-1} \mathcal{A}_i(\omega))) + \lambda_M \Pr(n|\omega) \log(\Pr(n|\omega)) \Big).$$

The equality of (11) and (12) follows from the symmetry of mutual information, defined in [Appendix 1](#). ■

I now resume the proof of [Lemma 6](#). First notice that [Lemma 14](#) establishes $\mathbf{C}(\mathbb{P}, \mu)$. For each $n \in \mathcal{N}$, let s_n denote a signal in \mathcal{S}_n which results in the posterior generated by signals in \mathcal{S}_n with probability one (in [Lemma 13](#) I show I can do this). Then notice:

$$\sum_{\omega \in \Omega} \int_s V(s) F(ds|\omega) \mu(\omega) = \sum_{n \in \mathcal{N}} V(s_n) \int_{s \in \mathcal{S}_n} \sum_{\omega \in \Omega} F(ds|\omega) \mu(\omega)$$

$$= \sum_{n \in \mathcal{N}} V(s_n) \Pr(n) = \sum_{n \in \mathcal{N}} \sum_{\omega \in \Omega} \mathbf{v}_n(\omega) F(\omega|s_n) \Pr(n)$$

$$= \sum_{n \in \mathcal{N}} \sum_{\omega \in \Omega} \mathbf{v}_n(\omega) \Pr(n|\omega) \mu(\omega)$$

Where the last step follows from the fact that $\Pr(X|Y)\Pr(Y) = \Pr(Y|X)\Pr(X)$. I now proceed with two proofs by contradiction. First, assume that (F, a) is a solution to (2) subject to (3),

which achieves expected utility U_1 , and let \mathbb{P} be the choice probabilities induced by it. Assume that \mathbb{P} is not a solution to (7) subject to (8) and (9), and thus there is a $\tilde{\mathbb{P}}$ which satisfies (8) and (9) and achieves expected utility $U_2 > U_1$. However, a strategy pairing (\tilde{F}, \tilde{a}) can be created that generates $\tilde{\mathbb{P}}$. For instance, for each of N distinct signals s_n , let $\tilde{a}(\tilde{F}(\omega|s_n)) \equiv n$, and let $\tilde{F}(s_n, \omega) = \tilde{\text{Pr}}(n|\omega)\mu(\omega) \forall \omega$ so that (3) is satisfied. This is impossible though as then (\tilde{F}, \tilde{a}) achieves $U_2 > U_1$ and (F, a) cannot have been optimal.

Similarly, assume that \mathbb{P} is a solution to (7) subject to (8) and (9), which achieves expected utility U_3 , but is not induced by a solution to (2) subject to (3). That is there is a \tilde{F} which satisfies (3) and achieves $U_4 > U_3$. This means, however, that $\tilde{\text{Pr}}(n|\omega) = \frac{\tilde{F}(s_n, \omega)}{\mu(\omega)}$ also achieves U_4 , which is impossible as \mathbb{P} was supposedly optimal and $\tilde{\mathbb{P}}$ satisfies (8) and (9). ■

Proof of Theorem 2. The Lagrangian for the above problem can be written:

$$\begin{aligned} \mathcal{L} = & \sum_{n \in \mathcal{N}} \sum_{\omega \in \Omega} \mathbf{v}_n(\omega) \text{Pr}(n|\omega) \mu(\omega) - \mathbf{C}(\mathbb{P}, \mu) + \sum_{n \in \mathcal{N}} \sum_{\omega \in \Omega} \xi_n(\omega) \text{Pr}(n|\omega) \mu(\omega) \\ & - \sum_{\omega \in \Omega} \gamma(\omega) \left(\sum_{n \in \mathcal{N}} \text{Pr}(n|\omega) - 1 \right) \mu(\omega). \end{aligned}$$

$\xi_n(\omega) \geq 0$ are the multipliers for (8), and $\gamma(\omega)$ are the multipliers for (9). If $\text{Pr}(n) = 0$, then $\text{Pr}(n|\omega) = 0 \forall \omega \in \Omega$. If $\text{Pr}(n|\cap_{i=1}^m \mathcal{A}_i(\omega)) = 0$ for some $m \in \{1, \dots, M-1\}$ and ω , then $\text{Pr}(n|\omega) = 0$. If $\text{Pr}(n) > 0$, and $\text{Pr}(n|\cap_{i=1}^m \mathcal{A}_i(\omega)) > 0, \forall m \in \{1, \dots, M-1\}$, then the first order condition with respect to $\text{Pr}(n|\omega)$ implies:

$$\begin{aligned} & \mathbf{v}_n(\omega) + \lambda_1(1 + \log \text{Pr}(n)) + (\lambda_2 - \lambda_1)(1 + \log \text{Pr}(n|\mathcal{A}_1(\omega))) \\ & + \dots + (\lambda_M - \lambda_{M-1})(1 + \log \text{Pr}(n|\cap_{i=1}^{M-1} \mathcal{A}_i(\omega))) - \lambda_M(1 + \log \text{Pr}(n|\omega)) = \gamma(\omega) - \xi_n(\omega) \end{aligned}$$

which then implies $\text{Pr}(n|\omega) > 0$ and $\xi_n(\omega) = 0$, because if not, and $\text{Pr}(n|\omega) = 0$, then since $\xi_n(\omega) \geq 0$, equality of the first order condition then necessitates $\gamma(\omega) = \infty$. This is impossible, however, since then $\forall \nu \in \mathcal{N}$ their respective first order conditions holding necessitates $\text{Pr}(\nu|\omega) = 0$. This being true $\forall \nu \in \mathcal{N}$ of course then violates (9). Thus, the first order condition implies:

$$\text{Pr}(n|\omega) = \text{Pr}(n)^{\frac{\lambda_1}{\lambda_M}} \text{Pr}(n|\mathcal{A}_1(\omega))^{\frac{\lambda_2 - \lambda_1}{\lambda_M}} \dots \text{Pr}(n|\cap_{i=1}^{M-1} \mathcal{A}_i(\omega))^{\frac{\lambda_M - \lambda_{M-1}}{\lambda_M}} e^{\frac{\mathbf{v}_n(\omega)}{\lambda_M}} e^{\frac{-\gamma(\omega)}{\lambda_M}} \quad (13)$$

Plugging (13) into (9), one can solve for $\gamma(\omega)$. Plugging $\gamma(\omega)$ back into (13) achieves the desired

result. ■

Proof of Corollary 1. Plug equation (10) into equation (7). ■

Proof of Theorem 3. Assume behavior is consistent with Theorem 2. The Lagrangian for the problem described in Corollary 1 is:

$$\begin{aligned} \mathcal{L} = & \sum_{\omega \in \Omega} \left(\log \left(\sum_{n=1}^N \Pr(n)^{\frac{\lambda_1}{\lambda_M}} \Pr(n | \mathcal{P}_{\lambda_1}(\omega))^{\frac{\lambda_2 - \lambda_1}{\lambda_M}} \dots \Pr(n | \cap_{i=1}^{M-1} \mathcal{P}_{\lambda_i}(\omega))^{\frac{\lambda_M - \lambda_{M-1}}{\lambda_M}} e^{\frac{v_n(\omega)}{\lambda_M}} \right) \mu(\omega) \right) \\ & + \sum_{A \in \times \{\mathcal{P}_{\lambda_i}\}_{i=1}^{M-1}} \sum_{n \in \mathcal{N}} \xi_n(A) \Pr(n|A) - \sum_{A \in \times \{\mathcal{P}_{\lambda_i}\}_{i=1}^{M-1}} \gamma(A) \left(\sum_{n \in \mathcal{N}} \Pr(n|A) - 1 \right) \end{aligned}$$

Using Theorem 2, the first order condition with respect to $\Pr(n|A)$ for some $A \in \times \{\mathcal{P}_{\lambda_i}\}_{i=1}^{M-1}$ and $\tilde{\omega} \in A$ is then:

$$\begin{aligned} & \left(\sum_{\omega \in \Omega} \frac{\lambda_1 \mu(\cap_{i=1}^{M-1} \mathcal{P}_{\lambda_i}(\tilde{\omega}))}{\lambda_M \Pr(n)} \Pr(n|\omega) \mu(\omega) \right) + \left(\sum_{\omega \in \mathcal{P}_{\lambda_1}(\tilde{\omega})} \frac{(\lambda_2 - \lambda_1) \mu(\cap_{i=1}^{M-1} \mathcal{P}_{\lambda_i}(\tilde{\omega}))}{\lambda_M \Pr(n | \mathcal{P}_{\lambda_1}(\omega)) \mu(\mathcal{P}_{\lambda_1}(\tilde{\omega}))} \Pr(n|\omega) \mu(\omega) \right) \\ & + \dots + \left(\sum_{\omega \in \cap_{i=1}^{M-1} \mathcal{P}_{\lambda_i}(\tilde{\omega})} \frac{(\lambda_M - \lambda_{M-1}) \mu(\cap_{i=1}^{M-1} \mathcal{P}_{\lambda_i}(\tilde{\omega}))}{\lambda_M \Pr(n | \cap_{i=1}^{M-1} \mathcal{P}_{\lambda_i}(\omega)) \mu(\cap_{i=1}^{M-1} \mathcal{P}_{\lambda_i}(\tilde{\omega}))} \Pr(n|\omega) \mu(\omega) \right) + \xi_n(A) = \gamma(A) \end{aligned}$$

For $n \in \mathcal{N}$ with $\Pr(n|A) > 0$, which there must be at least one of, $\xi_n(A) = 0$, and:

$$\sum_{\omega \in \Omega} \frac{\lambda_1 \mu(\cap_{i=1}^{M-1} \mathcal{P}_{\lambda_i}(\tilde{\omega}))}{\lambda_M \Pr(n)} \Pr(n|\omega) \mu(\omega) = \frac{\lambda_1 \mu(A)}{\lambda_M},$$

and for for each $m \in \{1, \dots, M-1\}$:

$$\sum_{\omega \in \cap_{i=1}^m \mathcal{P}_{\lambda_i}(\tilde{\omega})} \frac{(\lambda_{m+1} - \lambda_m) \mu(\cap_{i=1}^{M-1} \mathcal{P}_{\lambda_i}(\tilde{\omega}))}{\lambda_M \Pr(n | \cap_{i=1}^m \mathcal{P}_{\lambda_i}(\omega)) \mu(\cap_{i=1}^m \mathcal{P}_{\lambda_i}(\tilde{\omega}))} \Pr(n|\omega) \mu(\omega) = \frac{(\lambda_{m+1} - \lambda_m) \mu(A)}{\lambda_M}.$$

This tells us:

$$\gamma(A) = \mu(A) \left(\frac{\lambda_1}{\lambda_M} + \frac{\lambda_2 - \lambda_1}{\lambda_M} + \dots + \frac{\lambda_M - \lambda_{M-1}}{\lambda_M} \right) = \mu(A), \quad \forall A \in \times \{\mathcal{P}_{\lambda_i}\}_{i=1}^{M-1}.$$

I do, however, need to worry about options n that such that there is an event A with $\Pr(n|A) = 0$. I next show with a proof by contradiction that for all $n \in \mathcal{N}$ such that $\Pr(n) > 0$, and for all $A \in \times \{\mathcal{P}_{\lambda_i}\}_{i=1}^{M-1}$: $\Pr(n|A) > 0$ (which implies $\Pr(n|\tilde{\omega}) > 0$ for all $\tilde{\omega} \in \Omega$). Assume there is an alternative $n \in \mathcal{N}$ such that $\Pr(n) > 0$, and $\exists A \in \times \{\mathcal{P}_{\lambda_i}\}_{i=1}^{M-1}$ such that $\Pr(n|A) = 0$.

This means for some $\tilde{\omega} \in A$ there is a $m \in \{1, \dots, M-1\}$ such that $\Pr(n | \cap_{i=1}^m \mathcal{P}_{\lambda_i}(\tilde{\omega})) = 0$ and $\Pr(n | \cap_{i=1}^{m-1} \mathcal{P}_{\lambda_i}(\tilde{\omega})) > 0$. In this case the Karush-Khun-Tucker conditions are necessary, but they are not always sufficient since at corner solutions our objective is not differentiable (Lange, 2013). When considering the validity of a corner where $\Pr(n|A) = 0$, I must bound away from 0 and let the difference go to zero to ensure differentiability. $\forall B \in \times\{\mathcal{P}_{\lambda_i}\}_{i=1}^{M-1}$, if $\Pr(n|B) = 0$, instead let $\tilde{\Pr}(n|B) = \epsilon > 0$ (so that the multiplier $\xi_n(A)$ corresponds to the constraint $\Pr(n|A) \geq \epsilon$), and for each such B and for each of the k alternatives $\nu \in \mathcal{N} \setminus n$ such that $\Pr(\nu|B) > 0$, let $\tilde{\Pr}(\nu|B) = \Pr(\nu|B) - \epsilon/k > 0$ (so that probabilities sum to one), keeping all other probabilities the same. Then, the first order condition with respect to $\Pr(n|A)$ is thus:

$$\begin{aligned} & \left(\sum_{\omega \in \Omega} \frac{\lambda_1 \mu(\cap_{i=1}^{M-1} \mathcal{P}_{\lambda_i}(\tilde{\omega}))}{\lambda_M \tilde{\Pr}(n)} \tilde{\Pr}(n|\omega) \mu(\omega) \right) + \left(\sum_{\omega \in \mathcal{P}_{\lambda_1}(\tilde{\omega})} \frac{(\lambda_2 - \lambda_1) \mu(\cap_{i=1}^{M-1} \mathcal{P}_{\lambda_i}(\tilde{\omega}))}{\lambda_M \tilde{\Pr}(n|\mathcal{P}_{\lambda_1}(\omega)) \mu(\mathcal{P}_{\lambda_1}(\tilde{\omega}))} \tilde{\Pr}(n|\omega) \mu(\omega) \right) \\ & + \dots + \left(\sum_{\omega \in \cap_{i=1}^{M-1} \mathcal{P}_{\lambda_i}(\tilde{\omega})} \frac{(\lambda_M - \lambda_{M-1}) \mu(\cap_{i=1}^{M-1} \mathcal{P}_{\lambda_i}(\tilde{\omega}))}{\lambda_M \tilde{\Pr}(n | \cap_{i=1}^{M-1} \mathcal{P}_{\lambda_i}(\omega)) \mu(\cap_{i=1}^{M-1} \mathcal{P}_{\lambda_i}(\tilde{\omega}))} \tilde{\Pr}(n|\omega) \mu(\omega) \right) \\ & + \xi_n(A) = \mu(A), \end{aligned}$$

which cannot be satisfied for $\tilde{\Pr}(n|A) = \epsilon$ since $\xi_n(A) \geq 0$, and the last sum (at the very least) goes to infinite as ϵ goes to zero. Thus, the behavior is optimal when it is consistent with Theorem 2 as long as the consideration set is optimal, since then I can ignore the options that are not considered and differentiability implies the Karush-Khun-Tucker conditions are necessary and sufficient (Lange, 2013).

Finally, I need to consider the validity of a corner where $\Pr(\nu) = 0$ (so I can evaluate if the consideration set is optimal). I must again bound away from $\Pr(\nu) = 0$ and let the difference go to zero to ensure differentiability so the Karush-Khun-Tucker conditions are necessary and sufficient (Lange, 2013). For each of the r options $\nu \in \mathcal{N}$ such that $\Pr(\nu) = 0$, $\forall B \in \times\{\mathcal{P}_{\lambda_i}\}_{i=1}^{M-1}$, let $\tilde{\Pr}(\nu|B) = \frac{\epsilon}{r} > 0$ (so that the multiplier $\xi_n(A)$ corresponds to the constraint $\Pr(\nu|A) \geq \frac{\epsilon}{r}$), and for each such B and for each of the k alternatives $n \in \mathcal{N}$ such that $\Pr(n) > 0$, let $\tilde{\Pr}(n|B) = \Pr(n|B) - \epsilon/k > 0$ (so that probabilities sum to one). Picking ϵ small enough, the first order condition with respect to $\Pr(\nu|B)$ (ν such that $\Pr(\nu) = 0$) gives:

$$\begin{aligned}
& \frac{\lambda_1}{\lambda_M} \left(\sum_{\omega \in \Omega} \frac{e^{\frac{\mathbf{v}_\nu(\omega)}{\lambda_M}} \mu(B)}{\sum_{n \in \mathcal{N}} \tilde{\text{Pr}}(n)^{\frac{\lambda_1}{\lambda_M}} (\tilde{\text{Pr}}(n | \mathcal{P}_{\lambda_1}(\omega)))^{\frac{\lambda_2 - \lambda_1}{\lambda_M}} \dots (\tilde{\text{Pr}}(n | \cap_{i=1}^{M-1} \mathcal{P}_{\lambda_i}(\omega)))^{\frac{\lambda_M - \lambda_{M-1}}{\lambda_M}} e^{\frac{\mathbf{v}_n(\omega)}{\lambda_M}}} \mu(\omega) \right) \\
& + \frac{\lambda_2 - \lambda_1}{\lambda_M} \left(\sum_{\omega \in \mathcal{P}_{\lambda_1}(\tilde{\omega})} \frac{e^{\frac{\mathbf{v}_\nu(\omega)}{\lambda_M}} \mu(B)}{\sum_{n \in \mathcal{N}} \tilde{\text{Pr}}(n)^{\frac{\lambda_1}{\lambda_M}} (\tilde{\text{Pr}}(n | \mathcal{P}_{\lambda_1}(\omega)))^{\frac{\lambda_2 - \lambda_1}{\lambda_M}} \dots (\tilde{\text{Pr}}(n | \cap_{i=1}^{M-1} \mathcal{P}_{\lambda_i}(\omega)))^{\frac{\lambda_M - \lambda_{M-1}}{\lambda_M}} e^{\frac{\mathbf{v}_n(\omega)}{\lambda_M}}} \right. \\
& \quad \left. \cdot \mu(\omega | \mathcal{P}_{\lambda_1}(\tilde{\omega})) \right) \\
& + \dots + \frac{\lambda_M - \lambda_{M-1}}{\lambda_M} \left(\sum_{\omega \in B} \frac{e^{\frac{\mathbf{v}_\nu(\omega)}{\lambda_M}} \mu(B)}{\sum_{n \in \mathcal{N}} \tilde{\text{Pr}}(n)^{\frac{\lambda_1}{\lambda_M}} (\tilde{\text{Pr}}(n | \mathcal{P}_{\lambda_1}(\omega)))^{\frac{\lambda_2 - \lambda_1}{\lambda_M}} \dots (\tilde{\text{Pr}}(n | \cap_{i=1}^{M-1} \mathcal{P}_{\lambda_i}(\omega)))^{\frac{\lambda_M - \lambda_{M-1}}{\lambda_M}} e^{\frac{\mathbf{v}_n(\omega)}{\lambda_M}}} \right. \\
& \quad \left. \cdot \mu(\omega | B) \right) \\
& + \xi_\nu(B) = \mu(B).
\end{aligned}$$

This establishes the desired results when ϵ goes to zero since $\xi_\nu(B) \geq 0$. ■

Proof of Theorem 4. A fixed effect interpretation of MASE follows easily from the optimal choice probabilities described in Theorem 2:

$$\begin{aligned}
\text{Pr}(n|\omega) &= \frac{\text{Pr}(n)^{\frac{\lambda_1}{\lambda_M}} \text{Pr}(n|\mathcal{A}_1(\omega))^{\frac{\lambda_2 - \lambda_1}{\lambda_M}} \dots \text{Pr}(n|\cap_{i=1}^{M-1} \mathcal{A}_i(\omega))^{\frac{\lambda_M - \lambda_{M-1}}{\lambda_M}} e^{\frac{\mathbf{v}_n(\omega)}{\lambda_M}}}{\sum_{\nu \in \mathcal{N}} \text{Pr}(\nu)^{\frac{\lambda_1}{\lambda_M}} \text{Pr}(\nu|\mathcal{A}_1(\omega))^{\frac{\lambda_2 - \lambda_1}{\lambda_M}} \dots \text{Pr}(\nu|\cap_{i=1}^{M-1} \mathcal{A}_i(\omega))^{\frac{\lambda_M - \lambda_{M-1}}{\lambda_M}} e^{\frac{\mathbf{v}_\nu(\omega)}{\lambda_M}}} \\
&= \frac{(\text{NPr}(n))^{\frac{\lambda_1}{\lambda_M}} (\text{NPr}(n|\mathcal{A}_1(\omega)))^{\frac{\lambda_2 - \lambda_1}{\lambda_M}} \dots (\text{NPr}(n|\cap_{i=1}^{M-1} \mathcal{A}_i(\omega)))^{\frac{\lambda_M - \lambda_{M-1}}{\lambda_M}} e^{\frac{\mathbf{v}_n(\omega)}{\lambda_M}}}{\sum_{\nu \in \mathcal{N}} (\text{NPr}(\nu))^{\frac{\lambda_1}{\lambda_M}} (\text{NPr}(\nu|\mathcal{A}_1(\omega)))^{\frac{\lambda_2 - \lambda_1}{\lambda_M}} \dots (\text{NPr}(\nu|\cap_{i=1}^{M-1} \mathcal{A}_i(\omega)))^{\frac{\lambda_M - \lambda_{M-1}}{\lambda_M}} e^{\frac{\mathbf{v}_\nu(\omega)}{\lambda_M}}} \\
&= \frac{e^{\frac{\mathbf{v}_n(\omega) + \lambda_1 \alpha_n^0 + (\lambda_2 - \lambda_1) \alpha_n^1 + \dots + (\lambda_M - \lambda_{M-1}) \alpha_n^{M-1}}{\lambda_M}}}{\sum_{\nu \in \mathcal{N}} e^{\frac{\mathbf{v}_\nu(\omega) + \lambda_1 \alpha_\nu^0 + (\lambda_2 - \lambda_1) \alpha_\nu^1 + \dots + (\lambda_M - \lambda_{M-1}) \alpha_\nu^{M-1}}{\lambda_M}}}
\end{aligned}$$

Where $\alpha_\nu^0 = \log(\text{NPr}(\nu))$, and for $m \in \{1, \dots, M-1\}$ I have $\alpha_\nu^m = \log(\text{NPr}(\nu | \cap_{i=1}^m \mathcal{A}_i(\omega)))$. Normalizing the value of the options by λ_M , namely letting $\tilde{v}_n = \frac{\mathbf{v}_n(\omega)}{\lambda_M}$, and defining α_n appropriately,

agent choice behavior described by RI with MASE can then be interpreted as a RU model where each option n has perceived value:

$$u_n = \tilde{v}_n + \frac{\lambda_1}{\lambda_M} \alpha_n^0 + \frac{\lambda_2 - \lambda_1}{\lambda_M} \alpha_n^1 + \dots + \frac{\lambda_M - \lambda_{M-1}}{\lambda_M} \alpha_n^{M-1} + \epsilon_n = \tilde{v}_n + \alpha_n + \epsilon_n$$

The only kind of RU model consistent with this behavior is one where ϵ_n is distributed iid according to a Gumbel distribution (Train, 2009). ■

Appendix 3

An Introduction to Rational Inattention with Shannon Entropy

In the rational inattention (RI) literature learning by the agent is typically modelled as the choice of a signal structure, which means the agent chooses the probability of receiving different signals in different states of the world. Receiving a signal updates the agent's belief about the state of the world, giving them a more informed posterior belief. More informative signal structures are more costly for the agent, but allow them to make a more informed decision about which option to select.

Suppose that the uncertainty faced by the agent is described by a measurable space (Ω, \mathcal{F}) , where Ω is a finite set of possible **states of the world** (the state space), and \mathcal{F} is the set of **events** generated by Ω (the power set of Ω). I call $\mu : \mathcal{F} \rightarrow [0, 1]$, which assigns probabilities to events, the **prior** belief of the agent.

Suppose that an agent that has stopped learning must make a selection from a set of **options**, denoted $\mathcal{N} = \{1, \dots, N\}$. Each option, $n \in \mathcal{N}$, in each state of the world, $\omega \in \Omega$, has a (finite) **value** to the agent $\mathbf{v}_n(\omega)$.

The agent's problem is to maximize the expected value of their selected option less the cost of learning. They do this by choosing an **information strategy** $F(s, \omega) \in \Delta(\mathbb{R} \times \Omega)$, which is a joint distribution between s , the observed **signal**, and the states of the world.²² The only restriction on the information strategy is that the marginal, $F(\omega) : \mathcal{F} \rightarrow \mathbb{R}_+$, must equal the prior μ . Alternatively, an agent can select a probability measure $F(s|\omega) : \mathbb{R} \rightarrow \mathbb{R}_+$ for each $\omega \in \Omega$, which, combined with μ , determine both $F(s, \omega)$ and the posterior $F(\omega|s)$. It is a property of the

²²The decision to allow s to be any real number is rather arbitrary. This is a much richer signal space than is required in practice. It is shown that an optimal strategy only results in one of at most N different signals s being observed.

cost function for information derived in this paper, as is true with Shannon Entropy, that if $F(s, \omega)$ is optimal, then the agent is done learning after a single signal s . After the signal is realized, the agent simply picks the action with the highest expected value:

$$a(s|F) = \arg \max_{n \in \mathcal{N}} \mathbb{E}_{F(\omega|s)}[\mathbf{v}_n(\omega)].$$

Ignoring the cost of learning momentarily, the value to the agent of receiving a signal s , which induces posterior $F(\omega|s)$, is then:

$$V(s|F) = \max_{n \in \mathcal{N}} \mathbb{E}_{F(\omega|s)}[\mathbf{v}_n(\omega)].$$

Let the expected cost of a particular information strategy, given the agent's prior, be denoted $\mathbf{C}(F(s, \omega), \mu)$. I describe the form of the cost functions studied in this paper in [Section 3.1](#). The agent's problem can thus be written:

$$\max_{F \in \Delta(\mathbb{R} \times \Omega)} \sum_{\omega \in \Omega} \int_s V(s|F) F(ds|\omega) \mu(\omega) - \mathbf{C}(F(s, \omega), \mu),$$

$$\text{such that } \forall \omega \in \Omega : \int_s F(ds, \omega) = \mu(\omega).$$

The choice behavior the agent exhibits depends on the cost function for information. Shannon Entropy is a measure of uncertainty with an axiomatic foundation that can be used to assign costs to information and is frequently used in practice. If I am given a partition of the possible states of the world $\mathcal{P} = \{A_1, \dots, A_m\}$, and probability measure μ over these events, the uncertainty about which event has occurred, as measured by **Shannon Entropy**, is defined:²³

$$\mathcal{H}(\mathcal{P}, \mu) = - \sum_{i=1}^m \mu(A_i) \log(\mu(A_i)). \quad (14)$$

The convention used here is to set $0 \log(0) = 0$.

If an agent has prior μ about the state of the world, and their beliefs are updated to the posterior $\mu(\cdot|s)$ after they receive a signal s , then there is a change in the uncertainty as measured by Shannon Entropy. In the Shannon model of RI, the cost of an information strategy $F(s, \omega)$ is

²³This measure is only unique up to a positive multiplier.

measured as the expected reduction in total uncertainty as measured by Shannon Entropy:

$$\mathbb{E}\left[\mathcal{H}(\mathcal{P}, \mu) - \mathcal{H}(\mathcal{P}, \mu(\cdot|s))\right],$$

where $\mathcal{P} = \{\{\omega_1\}, \{\omega_2\}, \dots, \{\omega_n\}\}$ is the finest partition of the state space. Bayes rule, and the nature of Shannon Entropy, guarantee that every potential information strategy has a weakly positive cost. You can click [here](#) if you want to return to [Section 2](#).

A Discussion of Theorem 2

The behavior described in [Theorem 2](#) has many intuitive features. It is also a quite natural extension of the analogous result from [Matějka and McKay \(2015\)](#), which is described in equation (15). If I assume the agent has prior μ , and the cost of information is measured with Shannon Entropy so there is only one attribute to learn about (the environment studied in [Matějka and McKay \(2015\)](#)) that has associated multiplier λ_2 , then if the agent does optimal research in state $\omega \in \Omega$, they select option n from \mathcal{N} with probability:

$$\Pr(n|\omega) = \frac{\Pr(n)e^{\frac{v_n(\omega)}{\lambda_2}}}{\sum_{\nu \in \mathcal{N}} \Pr(\nu)e^{\frac{v_\nu(\omega)}{\lambda_2}}}. \quad (15)$$

One takeaway from the formula in (15) is that when Shannon Entropy is used to measure uncertainty the chance of the agent selecting an option n in a particular state of the world ω is fully determined by the unconditional chances of the options being selected, $\Pr(n)$, and the realized values of the options in that state of the world. Beyond this takeaway, the formula in (15) also has many intuitive features. If λ_2 grows (shrinks), which represents an increase (decrease) in the difficulty of learning, the value of each option in the realized state becomes less (more) significant for the determination of the selected option, and the significance of the agent's prior increases (decreases). If λ_2 approaches infinity, the realized values become insignificant, and the behavior of the agent approaches the behavior of the agent in the case where learning is impossible: they choose their option based on their prior. If λ_2 approaches zero the unconditional priors become insignificant, and the behavior of the agent approaches the behavior of the agent in the case where learning is costless: they choose the option with the highest realized value.

If I instead assume that the cost of information is measured with MASE and the agent may also learn about another attribute \mathcal{A}_1 with a lower associated multiplier λ_1 , then if $\mathcal{A}_1 \neq \Omega$, and

the agent does optimal research in state $\omega \in \Omega$, they select option n from their set of options \mathcal{N} with probability:

$$\Pr(n|\omega) = \frac{\Pr(n)^{\frac{\lambda_1}{\lambda_2}} \Pr(n|\mathcal{A}_1(\omega))^{\frac{\lambda_2-\lambda_1}{\lambda_2}} e^{\frac{\mathbf{v}_n(\omega)}{\lambda_2}}}{\sum_{\nu \in \mathcal{N}} \Pr(\nu)^{\frac{\lambda_1}{\lambda_2}} \Pr(\nu|\mathcal{A}_1(\omega))^{\frac{\lambda_2-\lambda_1}{\lambda_2}} e^{\frac{\mathbf{v}_\nu(\omega)}{\lambda_2}}}. \quad (16)$$

With MASE, as the formula in (16) indicates, the chance of the agent selecting an option n in a particular state of the world ω depends not only on the unconditional chances of the options being selected and the realized values of the options, but also on the realized value of \mathcal{A}_1 . When option n is in general desirable in $\mathcal{A}_1(\omega)$ relative to the other options, then $\Pr(n|\mathcal{A}_1(\omega))$ is larger, and there may be a high chance of n being selected, even if $\Pr(n)$ is not that large, and $\mathbf{v}_n(\omega)$ is not that high.

The formula in (16) also has many intuitive features. It maintains the intuitive comparative statistics for λ_2 that the formula in (15) had, and also features intuitive properties for $\Pr(n|\mathcal{A}_1(\omega))$ and λ_1 .

If observing $\mathcal{A}_1(\omega)$ is completely uninformative about the value of the options, then it is optimal for the agent to select $\Pr(n|\mathcal{A}_1(\omega)) = \Pr(n)$ since $\min_{S \in \mathcal{S}^b(\Omega)} C(S, \mu)$ is strictly concave in μ . In this case $\Pr(n)^{\frac{\lambda_1}{\lambda_2}} \Pr(n|\mathcal{A}_1(\omega))^{\frac{\lambda_2-\lambda_1}{\lambda_2}} = \Pr(n)$, and behavior is identical to that in (15). If the cheaper to learn about attribute is irrelevant information it is thus ignored, and behavior collapses back to the environment described in Matějka and McKay (2015), as I should desire.

If λ_1 approaches λ_2 (the cheaper to learn about attributes becomes close to as expensive as the more expensive to learn about attributes) then behavior approaches that described in (15) since $\Pr(n)^{\frac{\lambda_1}{\lambda_2}} \Pr(n|\mathcal{A}_1(\omega))^{\frac{\lambda_2-\lambda_1}{\lambda_2}} \rightarrow \Pr(n)$. Thus, if an insignificantly cheaper to learn about attribute is introduced behavior is changed in an insignificant fashion (see Figure 1). Again, this seems like a desirable property.

If λ_1 approaches zero then the role of the unconditional priors dissipates, and exponent on $\Pr(n|\mathcal{A}_1(\omega))$ approaches one, meaning it replaces the unconditional prior from (15). This makes sense because if λ_1 goes to zero it means $\mathcal{A}_1(\omega)$ can essentially be viewed for free, in which case behavior within each $\mathcal{A}_1(\omega)$ should resemble that in the setting where there is only one attribute with multiplier λ_2 and a prior of $\mu(\cdot|\mathcal{A}_1(\omega))$.

I can add new partitions with new multipliers and the description of behavior in Theorem 2 maintains the intuitive properties described in the paragraphs above. RI with MASE is thus a very natural extension of RI with Shannon Entropy. You can click [here](#) to return to Theorem 2.