

Article

Entropic Updating of Probability and Density Matrices

Kevin Vanslette

Department of Physics, University at Albany (SUNY), Albany, NY 12222, USA

* Correspondence: kvanslette@albany.edu

Abstract: We find that the standard relative entropy and the Umegaki entropy are designed for the purpose of inferentially updating probability and density matrices respectively. From the same set of inferentially guided design criteria, both of the previously stated entropies are derived in parallel. This formulates a quantum maximum entropy method for the purpose of inferring density matrices in the absence of complete information.

Keywords: probability theory; entropy; quantum relative entropy; quantum information; quantum mechanics; inference

1. Introduction

We design an inferential updating procedure for probability distributions and density matrices such that inductive inferences may be made. The inferential updating tools found in this derivation take the form of the standard and quantum relative entropy functionals, and thus we find the functionals are *designed* for the purpose of updating probability distributions and density matrices respectively. Design derivations which found the entropy to be a tool for inference originally required five *design criteria* (DC) [1–3], this was reduced to four in [4–6], and then down to three in [7]. We reduced the number of required DC down to two while also providing the first *design* derivation of the quantum relative entropy – *using the same design criteria and inferential principles in both instances*.

The designed quantum relative entropy takes the form of Umegaki's quantum relative entropy, and thus it has the "proper asymptotic form of the relative entropy in quantum (mechanics)" [8–10]. Recently, [11] gave an axiomatic characterization of the quantum relative entropy that "uniquely determines the quantum relative entropy". Our derivation differs from their's, again in that we *design* the quantum relative entropy for a purpose, but also that our DCs are imposed on what turns out to be the functional derivative of the quantum relative entropy rather than on the quantum relative entropy itself. The use of a quantum entropy for the purpose of inference has a large history: Jaynes [12,13] invented the notion of the quantum maximum entropy method [14], while it was perpetuated by [15–22] and many others. However, we find the quantum *relative* entropy to be the suitable entropy for updating density matrices, rather than the von Neumann. The relevant results of their papers may be found using our quantum relative entropy with a suitable uniform prior density matrix.

It should be noted that because the relative entropies were reached by design, they may be interpret as such, "the relative entropies are tools for updating", which means we no longer need to attach an interpretation *ex post facto* – as a measure of disorder or amount of missing information. In this sense, the relative entropies were built for the purpose of saturating their own interpretation [4,7].

The remainder of the paper is organized as follows: First we will discuss some universally applicable principles of inference and motivate the design of an entropy function able to rank probability distributions. This entropy function will be designed such it is consistent with inference by applying a few reasonable design criteria, which are guided by the aforementioned principles of inference. Using the same principles of inference and design criteria, we find the form of the quantum relative entropy suitable for inference. We end with concluding remarks.

38 Solutions for $\hat{\rho}$ by maximizing the quantum relative entropy give insight into the Quantum Bayes'
39 Rule in the sense of [23–26]. This, and a few other applications of the quantum maximum entropy
40 method, will be discussed in a future article.

41 2. The Design of Entropic Inference

42 Inference is the appropriate updating of probability distributions when new information is
43 received. Bayes' rule and Jeffrey's rule are both equipped to handle information in the form of data;
44 however, the updating of a probability distribution due to the knowledge of an expectation value was
45 realized by Jaynes [12–14] through the method of maximum entropy. The two methods for inference
46 were thought to be devoid of one another until the work of [27], which showed Bayes' and Jeffrey's
47 Rule to be consistent with the method of maximum entropy when the expectation values were in the
48 form of data [27]. In the spirit of the derivation we will carry-on as if the maximum entropy method
49 were not known and show how it may be derived as an application of inference.

50 Given a probability distribution $\varphi(x)$ over a general set of propositions $x \in X$, it is self evident
51 that if new information is learned, we are entitled to assign a new probability distribution $\rho(x)$ that
52 somehow reflects this new information while also respecting our prior probability distribution $\varphi(x)$.
53 The main question we must address is: "Given some information, to what posterior probability
54 distribution $\rho(x)$ should we update our prior probability distribution $\varphi(x)$ to?", that is,

$$\varphi(x) \xrightarrow{*} \rho(x)?$$

55 This specifies the problem of inductive inference. Since "information" has many colloquial,
56 yet potentially conflicting, definitions, we remove potential confusion by defining **information**
57 operationally (*) as the *rationale* that causes a probability distribution to change (inspired by and
58 adapted from [7]). Directly from [7]:

59
60 "Our goal is to design a method that allows a systematic search for the preferred posterior
61 distribution. The central idea, first proposed in [4] is disarmingly simple: to select the posterior first
62 rank all candidate distributions in increasing *order of preference* and then pick the distribution that
63 ranks the highest. Irrespective of what it is that makes one distribution preferable over another (we
64 will get to that soon enough) it is clear that any ranking according to preference must be transitive: if
65 distribution ρ_1 is preferred over distribution ρ_2 , and ρ_2 is preferred over ρ_3 , then ρ_1 is preferred over
66 ρ_3 . Such transitive rankings are implemented by assigning to each $\rho(x)$ a real number $S[\rho]$, which is
67 called the entropy of ρ , in such a way that if ρ_1 is preferred over ρ_2 , then $S[\rho_1] > S[\rho_2]$. The selected
68 distribution (one or possibly many, for there may be several equally preferred distributions) is that
69 which maximizes the entropy functional."

70
71 Because we wish to update from prior distributions φ to posterior distributions ρ by ranking, the
72 entropy functional $S[\rho, \varphi]$, is a real function of both φ and ρ . In the absence of new information, there
73 is no available *rationale* to prefer any ρ to the original φ , and thereby the relative entropy should be
74 designed such that the selected posterior is equal to the prior φ (in the absence of new information).
75 The prior information encoded in $\varphi(x)$ is valuable and we should not change it unless we are informed
76 otherwise. Due to our definition of information, and our desire for objectivity, we state the predominate
77 guiding principle for inductive inference:

78 The Principle of Minimal Updating (PMU):

79 *A probability distribution should only be updated to the extent required by the new information.*

80
81 This simple statement provides the foundation for inference [7]. If the updating of probability
82 distributions is to be done objectively, then possibilities should not be needlessly ruled out or

83 suppressed. Being informationally stingy, that we should only update probability distributions
 84 when the information requires it, pushes inductive inference toward objectivity. Thus using the PMU
 85 helps formulate a pragmatic (and objective) procedure for making inferences using (informationally)
 86 subjective probability distributions [28].

87 This method of inference is only as universal and general as its ability to apply *equally well* to
 88 *any* specific inference problem. The notion of "specificity" is the notion of statistical independence; a
 89 special case is only special in that it is separable from other special cases. The notion that systems may
 90 be "sufficiently independent" plays a central and deep-seated role in science and the idea that some
 91 things can be neglected and that not everything matters, is implemented by imposing criteria that tells
 92 us how to handle independent systems [7]. Ironically, the universally *shared* property by all specific
 93 inference problems is their ability to be *independent* of one another. Thus, a universal inference scheme
 94 based on the PMU permits,

95 Properties of Independence (PI):

96 *Subdomain Independence: When information is received about one set of propositions, it should not effect*
 97 *or change the state of knowledge (probability distribution) of the other propositions (else information was also*
 98 *received about them too);*

101 *And,*

99 *Subsystem Independence: When two systems are a-priori believed to be independent and we only receive*
 100 *information about one, then the state of knowledge of the other system remains unchanged.*

102 The PI's are special cases of the PMU that ultimately take the form of *design criteria* in the design
 103 derivation. The process of constraining the form of $S[\rho, \varphi]$ by imposing design criteria may be viewed
 104 as the process of *eliminative induction*, and after sufficient constraining, a single form for the entropy
 105 remains. Thus, the justification behind the surviving entropy is not that it leads to demonstrably
 106 correct inferences, but rather, that all other candidate entropies demonstrably fail to perform as desired
 107 [7]. Rather than the *design criteria* instructing one how to update, they instruct in what instances one
 108 should *not* update. That is, rather than justifying one way to skin a cat over another, we tell you when
 109 *not* to skin it, which is operationally unique – namely you don't do it – luckily enough for the cat.

110 2.1. The Design Criteria and the Standard Relative Entropy

111 The following *design criteria* (DC), guided by the PMU, are imposed and formulate the standard
 112 relative entropy as a tool for inference. The form of this presentation is inspired by [7].

113 DC1: Subdomain Independence

We keep the DC1 from [7] and review it below. DC1 imposes the first instance of when one should
 not update – the Subdomain PI. Suppose the information to be processed does *not* refer to a particular
 subdomain \mathcal{D} of the space \mathcal{X} of x 's. In the absence of new information about \mathcal{D} the PMU insists we do
 not change our minds about probabilities that are conditional on \mathcal{D} . Thus, we design the inference
 method so that $\varphi(x|\mathcal{D})$, the prior probability of x conditional on $x \in \mathcal{D}$, is not updated and therefore
 the selected conditional posterior is,

$$P(x|\mathcal{D}) = \varphi(x|\mathcal{D}). \quad (1)$$

114 (The notation will be as follows: we denote priors by φ , candidate posteriors by lower case ρ , and the
 115 selected posterior by upper case P .) We emphasize the point is not that we make the unwarranted
 116 assumption that keeping $\varphi(x|\mathcal{D})$ unchanged is guaranteed to lead to correct inferences. It need not;
 117 induction is risky. The point is, rather, that in the absence of any evidence to the contrary there is no
 118 reason to change our minds and the prior information takes priority.

119 DC1 Implementation:

120 Consider the set of microstates $x_i \in \mathcal{X}$ belonging to either of two non-overlapping domains \mathcal{D} or its
 121 compliment \mathcal{D}' , such that $\mathcal{X} = \mathcal{D} \cup \mathcal{D}'$ and $\emptyset = \mathcal{D} \cap \mathcal{D}'$. For convenience let $\rho(x_i) = \rho_i$. Consider the
 122 following constraints:

$$\rho(\mathcal{D}) = \sum_{i \in \mathcal{D}} \rho_i \quad \text{and} \quad \rho(\mathcal{D}') = \sum_{i \in \mathcal{D}'} \rho_i, \quad (2)$$

123 such that $\rho(\mathcal{D}) + \rho(\mathcal{D}') = 1$, and the following "local" constraints to \mathcal{D} and \mathcal{D}' respectively are,

$$\langle A \rangle = \sum_{i \in \mathcal{D}} \rho_i A_i \quad \text{and} \quad \langle A' \rangle = \sum_{i \in \mathcal{D}'} \rho_i A'_i. \quad (3)$$

124 As we are searching for the candidate distribution which maximizes S while obeying (2) and (3),
 125 we maximize the entropy $S \equiv S[\rho, \varphi]$ with respect to these expectation value constraints using the
 126 Lagrange multiplier method,

$$0 = \delta \left(S - \lambda [\rho(\mathcal{D}) - \sum_{i \in \mathcal{D}} \rho_i] - \mu [\langle A \rangle - \sum_{i \in \mathcal{D}} \rho_i A_i] \right. \\ \left. - \lambda' [\rho(\mathcal{D}') - \sum_{i \in \mathcal{D}'} \rho_i] - \mu' [\langle A' \rangle - \sum_{i \in \mathcal{D}'} \rho_i A'_i] \right),$$

127 and thus, the entropy is maximized when the following differential relationships hold:

$$\frac{\delta S}{\delta \rho_i} = \lambda + \mu A_i \quad \forall i \in \mathcal{D}, \quad (4)$$

$$\frac{\delta S}{\delta \rho_i} = \lambda' + \mu' A'_i \quad \forall i \in \mathcal{D}'. \quad (5)$$

128 Equations (2)-(5), are $n + 4$ equations we must solve to find the four Lagrange multipliers $\{\lambda, \lambda', \mu, \mu'\}$
 129 and the n probability values $\{\rho_i\}$.

130 If the subdomain constraint DC1 is imposed in the most restrictive case, then it will hold in
 131 general. The most restrictive case requires splitting \mathcal{X} into a set of $\{\mathcal{D}_i\}$ domains such that each \mathcal{D}_i
 132 singularly includes one microstate x_i . This gives,

$$\frac{\delta S}{\delta \rho_i} = \lambda_i + \mu_i A_i \quad \text{in each } \mathcal{D}_i. \quad (6)$$

133 Because the entropy $S = S[\rho_1, \rho_2, \dots; \varphi_1, \varphi_2, \dots]$ is a function over the probability of each microstate's
 134 posterior and prior distribution, its variational derivative is also a function of said probabilities in
 135 general,

$$\frac{\delta S}{\delta \rho_i} \equiv \phi_i(\rho_1, \rho_2, \dots; \varphi_1, \varphi_2, \dots) = \lambda_i + \mu_i A_i \quad \text{for each } (i, \mathcal{D}_i). \quad (7)$$

136 DC1 is imposed by constraining the form of $\phi_i(\rho_1, \rho_2, \dots; \varphi_1, \varphi_2, \dots) = \phi_i(\rho_i; \varphi_1, \varphi_2, \dots)$ to ensure that
 137 changes in $A_i \rightarrow A_i + \delta A_i$ have no influence over the value of ρ_j in domain \mathcal{D}_j , through ϕ_i , for $i \neq j$.
 138 If there is no new information about propositions in \mathcal{D}_j , its distribution should remain equal to φ_j
 139 by the PMU. We further restrict ϕ_i such that an arbitrary variation of $\varphi_j \rightarrow \varphi_j + \delta \varphi_j$ (a change in the
 140 prior state of knowledge of the microstate j) has no effect on ρ_i for $i \neq j$ and therefore DC1 imposes
 141 $\phi_i = \phi_i(\rho_i, \varphi_i)$, as is guided by the PMU. At this point it is easy to generalize the analysis to continuous
 142 microstates such that the indices become continuous $i \rightarrow x$, sums become integrals, and discrete
 143 probabilities become probability densities $\rho_i \rightarrow \rho(x)$.

144 **Remark:**

¹⁴⁵ We are designing the entropy for the purpose of ranking posterior probability distributions (for the
¹⁴⁶ purpose of inference); however, the highest ranked distribution is found by setting the variational
¹⁴⁷ derivative of $S[\rho, \varphi]$ equal to the variations of the expectation value constraints by the Lagrange
¹⁴⁸ multiplier method,

$$\frac{\delta S}{\delta \rho(x)} = \lambda + \sum_i \mu_i A_i(x). \quad (8)$$

¹⁴⁹ Therefore, the real quantity of interest is $\frac{\delta S}{\delta \rho(x)}$ rather than the specific form of $S[\rho, \varphi]$. All forms of
¹⁵⁰ $S[\rho, \varphi]$ that give the correct form of $\frac{\delta S}{\delta \rho(x)}$ are *equally valid* for the purpose of inference. Thus, every
¹⁵¹ design criteria may be made on the variational derivative of the entropy rather than the entropy itself,
¹⁵² which we do. When maximizing the entropy, for convenience, we will let,

$$\frac{\delta S}{\delta \rho(x)} \equiv \phi_x(\rho(x), \varphi(x)), \quad (9)$$

¹⁵³ and further use the shorthand $\phi_x(\rho, \varphi) \equiv \phi_x(\rho(x), \varphi(x))$, in all cases.

¹⁵⁴ **DC1':** *In the absence of new information, our new state of knowledge $\rho(x)$ is equal to the old state of knowledge*
¹⁵⁵ $\varphi(x)$.

¹⁵⁶ This is a special case of DC1, and is implemented differently than in [7]. The PMU is in principle a
¹⁵⁷ statement about informational honesty – that is, one should not “jump to conclusions” in light of new
¹⁵⁸ information and in the absence of new information, one should not change their state of knowledge.
¹⁵⁹ If no new information is given, the prior probability distribution $\varphi(x)$ does not change, that is, the
¹⁶⁰ posterior probability distribution $\rho(x) = \varphi(x)$ is equal to the prior probability. If we maximizing the
¹⁶¹ entropy without applying constraints,

$$\frac{\delta S}{\delta \rho(x)} = 0, \quad (10)$$

¹⁶² then DC1' imposes the following condition:

$$\frac{\delta S}{\delta \rho(x)} = \phi_x(\rho, \varphi) = \phi_x(\varphi, \varphi) = 0, \quad (11)$$

¹⁶³ for all x in this case. This special case of the DC1 and the PMU turns out to be incredibly constraining
¹⁶⁴ as we will see over the course of DC2.

¹⁶⁵ **Comment:**

¹⁶⁶ From [7]. If the variable x is continuous, DC1 requires that when information refers to points
¹⁶⁷ infinitely close but just outside the domain \mathcal{D} , that it will have no influence on probabilities conditional
¹⁶⁸ on \mathcal{D} . This may seem surprising as it may lead to updated probability distributions that are
¹⁶⁹ discontinuous. Is this a problem? No.

¹⁷⁰ In certain situations (e.g., physics) we might have explicit reasons to believe that conditions of
¹⁷¹ continuity or differentiability should be imposed and this information might be given to us in a variety
¹⁷² of ways. The crucial point, however – and this is a point that we keep and will keep reiterating – is
¹⁷³ that unless such information is explicitly given we should not assume it. If the new information leads
¹⁷⁴ to discontinuities, so be it.

¹⁷⁵

¹⁷⁶ **DC2: Subsystem Independence**

¹⁷⁷ DC2 imposes the second instance of when one should not update – the Subsystem PI. We
¹⁷⁸ emphasize that *DC2 is not a consistency requirement*. The argument we deploy is *not* that both the prior

179 and the new information tells us the systems are independent, in which case consistency requires that
 180 it should not matter whether the systems are treated jointly or separately. Rather, DC2 refers to a
 181 situation where the new information does not say whether the systems are independent or not, but
 182 information is given about each subsystem. The updating is being *designed* so that the independence
 183 reflected in the prior is maintained in the posterior by default via the PMU and the second clause of
 184 the PI's. [7]

185 The point is not that when we have no evidence for correlations we draw the firm conclusion that
 186 the systems must necessarily be independent. They could indeed have turned out to be correlated and
 187 then our inferences would be wrong. Again, induction involves risk. The point is rather that if the
 188 joint prior reflected independence and the new evidence is silent on the matter of correlations, then the
 189 prior takes precedence. As before, in this case subdomain independence, the probability distribution
 190 should not be updated unless the information requires it. [7]

191 **DC2 Implementation:**

192 Consider a composite system, $x = (x_1, x_2) \in \mathcal{X} = \mathcal{X}_1 \times \mathcal{X}_2$. Assume that all prior evidence led us
 193 to believe the subsystems are independent. This belief is reflected in the prior distribution: if the
 194 individual system priors are $\varphi_1(x_1)$ and $\varphi_2(x_2)$, then the prior for the whole system is their product
 195 $\varphi_1(x_1)\varphi_2(x_2)$. Further suppose that new information is acquired such that $\varphi_1(x_1)$ would by itself be
 196 updated to $P_1(x_1)$ and that $\varphi_2(x_2)$ would be itself be updated to $P_2(x_2)$. By design, the implementation
 197 of DC2 constrains the entropy functional such that in this case, the joint product prior $\varphi_1(x_1)\varphi_2(x_2)$
 198 updates to the selected product posterior $P_1(x_1)P_2(x_2)$. [7]

199 The argument below is considerably simplified if we expand the space of probabilities to include
 200 distributions that are not necessarily normalized. This does not represent any limitation because a
 201 normalization constraint may always be applied. We consider a few special cases below:

202

Case 1: We receive the extremely constraining information that the posterior distribution for system 1
 is completely specified to be $P_1(x_1)$ while we receive no information at all about system 2. We treat
 the two systems jointly. Maximize the joint entropy $S[\rho(x_1, x_2), \varphi(x_1)\varphi(x_2)]$ subject to the following
 constraints on the $\rho(x_1, x_2)$,

$$\int dx_2 \rho(x_1, x_2) = P_1(x_1). \quad (12)$$

Notice that the probability of each $x_1 \in \mathcal{X}_1$ within $\rho(x_1, x_2)$ is being constrained to $P_1(x_1)$ in the
 marginal. We therefore need a one Lagrange multiplier $\lambda_1(x_1)$ for each $x_1 \in \mathcal{X}_1$ to tie each value of
 $\int dx_2 \rho(x_1, x_2)$ to $P_1(x_1)$. Maximizing the entropy with respect to this constraint is,

$$\delta \left[S - \int dx_1 \lambda_1(x_1) \left(\int dx_2 \rho(x_1, x_2) - P_1(x_1) \right) \right] = 0, \quad (13)$$

which requires that

$$\lambda_1(x_1) = \phi_{x_1 x_2}(\rho(x_1, x_2), \varphi_1(x_1)\varphi_2(x_2)), \quad (14)$$

for arbitrary variations of $\rho(x_1, x_2)$. By design, DC2 is implemented by requiring $\varphi_1\varphi_2 \rightarrow P_1\varphi_2$ in this
 case, therefore,

$$\lambda_1(x_1) = \phi_{x_1 x_2}(P_1(x_1)\varphi_2(x_2), \varphi_1(x_1)\varphi_2(x_2)). \quad (15)$$

This equation must hold for all choices of x_2 and all choices of the prior $\varphi_2(x_2)$ as $\lambda_1(x_1)$ is independent
 of x_2 . Suppose we had chosen a different prior $\varphi'_2(x_2) = \varphi_2(x_2) + \delta\varphi_2(x_2)$ that disagrees with $\varphi_2(x_2)$.
 For all x_2 and $\delta\varphi_2(x_2)$, the multiplier $\lambda_1(x_1)$ remains unchanged as it constrains the independent
 $\rho(x_1) \rightarrow P_1(x_1)$. This means that any dependence that the right hand side might potentially have had
 on x_2 and on the prior $\varphi_2(x_2)$ must cancel out. This means that

$$\phi_{x_1 x_2}(P_1(x_1)\varphi_2(x_2), \varphi_1(x_1)\varphi_2(x_2)) = f_{x_1}(P_1(x_1), \varphi_1(x_1)). \quad (16)$$

Since φ_2 is arbitrary in f suppose further that we choose a constant prior set equal to one, $\varphi_2(x_2) = 1$, therefore

$$f_{x_1}(P_1(x_1), \varphi_1(x_1)) = \phi_{x_1 x_2}(P_1(x_1) * 1, \varphi_1(x_1) * 1) = \phi_{x_1}(P_1(x_1), \varphi_1(x_1)) \quad (17)$$

in general. This gives,

$$\lambda_1(x_1) = \phi_{x_1}(P_1(x_1), \varphi_1(x_1)). \quad (18)$$

203 The left hand side does not depend on x_2 , and therefore neither does the right hand side. An argument
204 exchanging systems 1 and 2 gives a similar result.

Case 1 - Conclusion: When the system 2 is not updated the dependence on φ_2 and x_2 drops out,

$$\phi_{x_1 x_2}(P_1(x_1) \varphi_2(x_2), \varphi_1(x_1) \varphi_2(x_2)) = \phi_{x_1}(P_1(x_1), \varphi_1(x_1)). \quad (19)$$

and vice-versa when system 1 is not updated,

$$\phi_{x_1 x_2}(\varphi_1(x_1) P_2(x_2), \varphi_1(x_1) \varphi_2(x_2)) = \phi_{x_2}(P_2(x_2), \varphi_2(x_2)). \quad (20)$$

205 As we seek the general functional form of $\phi_{x_1 x_2}$, and because the x_2 dependence drops out of (19)
206 and the x_1 dependence drops out of (20) for arbitrary φ_1, φ_2 and $\varphi_{12} = \varphi_1 \varphi_2$, the explicit coordinate
207 dependence in ϕ consequently drops out of both such that,

$$\phi_{x_1 x_2} \rightarrow \phi, \quad (21)$$

208 as $\phi = \phi(\rho(x), \varphi(x))$ must only depend on coordinates through the probability distributions
209 themselves. (As a double check, explicit coordinate dependence was included in the following
210 computations but inevitably dropped out due to the form the functional equations and DC1'. By the
211 argument above, and for simplicity, we drop the explicit coordinate dependence in ϕ here.)
212

Case 2: Now consider a different special case in which the marginal posterior distributions for systems 1 and 2 are both completely specified to be $P_1(x_1)$ and $P_2(x_2)$ respectively. Maximize the joint entropy $S[\rho(x_1, x_2), \varphi(x_1) \varphi(x_2)]$ subject to the following constraints on the $\rho(x_1, x_2)$,

$$\int dx_2 \rho(x_1, x_2) = P_1(x_1) \quad \text{and} \quad \int dx_1 \rho(x_1, x_2) = P_2(x_2). \quad (22)$$

213 Again, this is one constraint for each value of x_1 and one constraint for each value of x_2 , which therefore
214 require the separate multipliers $\mu_1(x_1)$ and $\mu_2(x_2)$. Maximizing S with respect to these constraints is
215 then,

$$\begin{aligned} 0 &= \delta \left[S - \int dx_1 \mu_1(x_1) \left(\int dx_2 \rho(x_1, x_2) - P_1(x_1) \right) \right. \\ &\quad \left. - \int dx_2 \mu_2(x_2) \left(\int dx_1 \rho(x_1, x_2) - P_2(x_2) \right) \right], \end{aligned} \quad (23)$$

216 leading to

$$\mu_1(x_1) + \mu_2(x_2) = \phi(\rho(x_1, x_2), \varphi_1(x_1) \varphi_2(x_2)). \quad (24)$$

The updating is being designed so that $\varphi_1 \varphi_2 \rightarrow P_1 P_2$, as the independent subsystems are being updated based on expectation values which are silent about correlations. DC2 thus imposes,

$$\mu_1(x_1) + \mu_2(x_2) = \phi(P_1(x_1) P_2(x_2), \varphi_1(x_1) \varphi_2(x_2)). \quad (25)$$

Write (25) as,

$$\mu_1(x_1) = \phi(P_1(x_1) P_2(x_2), \varphi_1(x_1) \varphi_2(x_2)) - \mu_2(x_2). \quad (26)$$

The left hand side is independent of x_2 so we can perform a trick similar to that we used before. Suppose we had chosen a different *constraint* $P'_2(x_2)$ that differs from $P_2(x_2)$ and a new prior $\varphi'_2(x_2)$ that differs from $\varphi_2(x_2)$ except at the value \bar{x}_2 . At the value \bar{x}_2 , the multiplier $\mu_1(x_1)$ remains unchanged for all $P'_2(x_2)$, $\varphi'_2(x_2)$, and thus x_2 . This means that any dependence that the right hand side might potentially have had on x_2 and on the choice of $P_2(x_2)$, $\varphi'_2(x_2)$ must cancel out leaving $\mu_1(x_1)$ unchanged. That is, the Lagrange multiplier $\mu(x_2)$ "pushes out" these dependences such that

$$\phi(P_1(x_1)P_2(x_2), \varphi_1(x_1)\varphi_2(x_2)) - \mu_2(x_2) = g(P_1(x_1), \varphi_1(x_1)). \quad (27)$$

Because $g(P_1(x_1), \varphi_1(x_1))$ is independent of arbitrary variations of $P_2(x_2)$ and $\varphi_2(x_2)$ on the LHS above – it is satisfied equally well for all choices. The form of $g = \phi(P_1(x_1), \varphi_1(x_1))$ is apparent if $P_2(x_2) = \varphi_2(x_2) = 1$ as $\mu_2(x_2) = 0$ similar to Case 1 as well as DC1'. Therefore, the Lagrange multiplier is

$$\mu_1(x_1) = \phi(P_1(x_1), \varphi_1(x_1)). \quad (28)$$

A similar analysis can be carried out for $\mu_2(x_2)$ leads to

$$\mu_2(x_2) = \phi(P_2(x_2), \varphi_2(x_2)). \quad (29)$$

Case 2 - Conclusion: Substituting back into (25) gives us a functional equation for ϕ ,

$$\phi(P_1P_2, \varphi_1\varphi_2) = \phi(P_1, \varphi_1) + \phi(P_2, \varphi_2). \quad (30)$$

²¹⁷ The general solution for this functional equation is derived in the Appendix, section 5.3, and is

$$\phi(\rho, \varphi) = a_1 \ln(\rho(x)) + a_2 \ln(\varphi(x)) \quad (31)$$

²¹⁸ where a_1, a_2 are constants. The constants are fixed by using DC1'. Letting $\rho_1(x_1) = \varphi_1(x_1) = \varphi_1$ gives
²¹⁹ $\phi(\varphi, \varphi) = 0$ by DC1', and therefore,

$$\phi(\varphi, \varphi) = (a_1 + a_2) \ln(\varphi) = 0, \quad (32)$$

²²⁰ so we are forced to conclude $a_1 = -a_2$ for arbitrary φ . Letting $a_1 \equiv A = -|A|$ such that we are really
²²¹ maximizing the entropy (although this is purely aesthetic) gives the general form of ϕ to be,

$$\phi(\rho, \varphi) = -|A| \ln\left(\frac{\rho(x)}{\varphi(x)}\right). \quad (33)$$

²²² As long as $A \neq 0$, the value of A is arbitrary as it always can be absorbed into the Lagrange multipliers.
²²³ The general form of the entropy designed for the purpose of inference of ρ is found by integrating ϕ ,
²²⁴ and therefore,

$$S(\rho(x), \varphi(x)) = -|A| \int dx (\rho(x) \ln\left(\frac{\rho(x)}{\varphi(x)}\right) - \rho(x)) + C[\varphi]. \quad (34)$$

²²⁵ The constant in ρ , $C[\varphi]$, will always drop out when varying ρ . The apparent extra term ($|A| \int \rho(x) dx$)
²²⁶ from integration cannot be dropped while simultaneously satisfying DC1', which requires $\rho(x) = \varphi(x)$
²²⁷ in the absence of constraints or when there is no change to one's information. In previous versions
²²⁸ where the integration term ($|A| \int \rho(x) dx$) is dropped, one obtains solutions like $\rho(x) = e^{-1} \varphi(x)$
²²⁹ (independent of whether $\varphi(x)$ was previously normalized or not) in the absence of new information.
²³⁰ Obviously this factor can be taken care of by normalization, and in this way both forms of the entropy
²³¹ are equally valid; however, this form of the entropy better adheres to the PMU through DC1'. Given

232 that we may regularly impose normalization, we may drop the extra $\int \rho(x)dx$ term and $C[\varphi]$. For
233 convenience then, (34) becomes

$$S(\rho(x), \varphi(x)) \rightarrow S^*(\rho(x), \varphi(x)) = -|A| \int dx \rho(x) \ln \left(\frac{\rho(x)}{\varphi(x)} \right), \quad (35)$$

234 which is a special case when the normalization constraint is being applied. Given normalization is
235 applied, the same selected posterior $\rho(x)$ maximizes both $S(\rho(x), \varphi(x))$ and $S^*(\rho(x), \varphi(x))$, and the
236 star notation may be dropped.

237 Remarks:

238 It can be seen that the relative entropy is invariant under coordinate transformations. This
239 implies that a system of coordinates carry no information and it is the "character" of the probability
240 distributions that are being ranked against one another rather than the specific set of propositions or
241 microstates they describe.

242 The general solution to the maximum entropy procedure with respect to N linear constraints in ρ ,
243 $\langle A_i(x) \rangle$, and normalization gives a canonical-like selected posterior probability distribution,

$$\rho(x) = \varphi(x) \exp \left(\sum_i \alpha_i A_i(x) \right). \quad (36)$$

244 The positive constant $|A|$ may always be absorbed into the Lagrange multipliers so we may let it equal
245 unity without loss of generality. DC1' is fully realized when we maximize with respect to a constraint
246 on $\rho(x)$ that is already held by $\varphi(x)$, such as $\langle x^2 \rangle = \int x^2 \rho(x)$ which happens to have the same value
247 as $\int x^2 \varphi(x)$, then its Lagrange multiplier is forcibly zero $\alpha_1 = 0$ (as can be seen in (36) using (34)), in
248 agreement with Jaynes. This gives the expected result $\rho(x) = \varphi(x)$ as there is no new information.
249 Our design has arrived at a refined maximum entropy method [12] as a universal probability updating
250 procedure [27].

251 3. The Design of the Quantum Relative Entropy

252 Last section we assumed that the universe of discourse (the set of relevant propositions or
253 microstates) $\mathcal{X} = \mathcal{A} \times \mathcal{B} \times \dots$ was known. In quantum physics things are a bit more ambiguous
254 because many probability distributions, or many experiments, can be associated to a given density
255 matrix. In this sense it helpful to think of density matrices as "placeholders" for probability distributions
256 rather than a probability distributions themselves. As any probability distribution from a given density
257 matrix, $\rho(\cdot) = \text{Tr}(|\cdot\rangle\langle\cdot| \hat{\rho})$, may be ranked using the standard relative entropy, it is unclear why we
258 would chose one universe of discourse over another. In lieu of this, such that one universe of discourse
259 is not given preferential treatment, we consider ranking entire density matrices against one another.
260 Probability distributions of interest may be found from the selected posterior density matrix. This
261 moves our universe of discourse from sets of propositions $\mathcal{X} \rightarrow \mathcal{H}$ to Hilbert space(s).

262 When the objects of study are quantum systems, we desire an objective procedure to update from
263 a prior density matrix $\hat{\rho}$ to a posterior density matrix $\hat{\rho}$. We will apply the same intuition for ranking
264 probability distributions (Section 2) and implement the PMU, PI, and design criteria to the ranking
265 of density matrices. We therefore find the quantum relative entropy $S(\hat{\rho}, \hat{\rho})$ to be designed for the
266 purpose of inferentially updating density matrices.

267 3.1. Designing the Quantum Relative Entropy

268 In this section we design the quantum relative entropy using the same inferentially guided *design*
269 *criteria* as were used in the standard relative entropy.

270 DC1: Subdomain Independence

271 The goal is to design a function $S(\hat{\rho}, \hat{\phi})$ which is able to rank density matrices. This insists that
 272 $S(\hat{\rho}, \hat{\phi})$ be a real scalar valued function of the posterior $\hat{\rho}$, and prior $\hat{\phi}$ density matrices, which we will
 273 call the quantum relative entropy or simply the entropy. An arbitrary variation of the entropy with
 274 respect to $\hat{\rho}$ is,

$$\delta S(\hat{\rho}, \hat{\phi}) = \sum_{ij} \frac{\delta S(\hat{\rho}, \hat{\phi})}{\delta \rho_{ij}} \delta \rho_{ij} = \sum_{ij} \left(\frac{\delta S(\hat{\rho}, \hat{\phi})}{\delta \hat{\rho}} \right)_{ij} \delta(\hat{\rho})_{ij} = \sum_{ij} \left(\frac{\delta S(\hat{\rho}, \hat{\phi})}{\delta \hat{\rho}^T} \right)_{ji} \delta(\hat{\rho})_{ij} = \text{Tr} \left(\frac{\delta S(\hat{\rho}, \hat{\phi})}{\delta \hat{\rho}^T} \delta \hat{\rho} \right). \quad (37)$$

275 We wish to maximize this entropy with respect to expectation value constraints, such as, $\langle A \rangle = \text{Tr}(\hat{A} \hat{\rho})$
 276 on $\hat{\rho}$. Using the Lagrange multiplier method to maximize the entropy with respect to $\langle A \rangle$ and
 277 normalization, is setting the variation equal to zero,

$$\delta \left(S(\hat{\rho}, \hat{\phi}) - \lambda [\text{Tr}(\hat{\rho}) - 1] - \alpha [\text{Tr}(\hat{A} \hat{\rho}) - \langle A \rangle] \right) = 0, \quad (38)$$

278 where λ and α are the Lagrange multipliers for the respective constraints. Because $S(\hat{\rho}, \hat{\phi})$ is a real
 279 number, we inevitably require δS to be real, but without imposing this directly, we find that requiring
 280 δS to be real requires $\hat{\rho}, \hat{A}$ to be Hermitian. At this point, it is simpler to allow for arbitrary variations
 281 of $\hat{\rho}$ such that,

$$\text{Tr} \left(\left(\frac{\delta S(\hat{\rho}, \hat{\phi})}{\delta \hat{\rho}^T} - \lambda \hat{1} - \alpha \hat{A} \right) \delta \hat{\rho} \right) = 0. \quad (39)$$

282 For these arbitrary variations, the variational derivative of S must satisfy,

$$\frac{\delta S(\hat{\rho}, \hat{\phi})}{\delta \hat{\rho}^T} = \lambda \hat{1} + \alpha \hat{A}, \quad (40)$$

283 at the maximum. As in the remark earlier, *all* forms of S which give the correct form of $\frac{\delta S(\hat{\rho}, \hat{\phi})}{\delta \hat{\rho}^T}$ under
 284 variation are *equally valid* for the purpose of inference. For notational convenience we let,

$$\frac{\delta S(\hat{\rho}, \hat{\phi})}{\delta \hat{\rho}^T} \equiv \phi(\hat{\rho}, \hat{\phi}), \quad (41)$$

285 which is a matrix valued function of the posterior and prior density matrices. The form of $\phi(\hat{\rho}, \hat{\phi})$ is
 286 already "local" in $\hat{\rho}$, so we don't need to constrain it further as we did in the original DC1.

287 **DC1':** *In the absence of new information, the new state $\hat{\rho}$ is equal to the old state $\hat{\phi}$.*

288 Applied to the ranking of density matrices, in the absence of new information, the density matrix
 289 $\hat{\phi}$ should not change, that is, the posterior density matrix $\hat{\rho} = \hat{\phi}$ is equal to the prior density matrix.
 290 Maximizing the entropy without applying any constraints gives,

$$\frac{\delta S(\hat{\rho}, \hat{\phi})}{\delta \hat{\rho}^T} = \hat{0}, \quad (42)$$

291 and therefore DC1' imposes the following condition in this case,

$$\frac{\delta S(\hat{\rho}, \hat{\phi})}{\delta \hat{\rho}^T} = \phi(\hat{\rho}, \hat{\phi}) = \phi(\hat{\phi}, \hat{\phi}) = \hat{0}. \quad (43)$$

292 As in the original DC1', if $\hat{\phi}$ is known to obey some expectation value constraint $\langle \hat{A} \rangle$, then if one goes
 293 out of their way to constrain $\hat{\rho}$ to that expectation value with nothing else, it follows from the PMU
 294 that $\hat{\rho} = \hat{\phi}$, as no information has been gained. This is not imposed directly, but can be verified later.

295 **DC2: Subsystem Independence**

296 The discussion of DC2 is the same as the standard relative entropy DC2 – it is not a consistency
 297 requirement, and the updating is *designed* so that the independence reflected in the prior is maintained
 298 in the posterior by default via the PMU, when the information provided is silent about correlations.

299 **DC2 Implementation:**

300 Consider a composite system living in the Hilbert space $\mathcal{H} = \mathcal{H}_1 \otimes \mathcal{H}_2$. Assume that all prior
 301 evidence led us to believe the systems were independent. This is reflected in the prior density matrix:
 302 if the individual system priors are $\hat{\rho}_1$ and $\hat{\rho}_2$, then the joint prior for the whole system is $\hat{\rho}_1 \otimes \hat{\rho}_2$.
 303 Further suppose that new information is acquired such that $\hat{\rho}_1$ would by itself be updated to $\hat{\rho}_1$ and
 304 that $\hat{\rho}_2$ would be itself be updated to $\hat{\rho}_2$. By design, the implementation of DC2 constrains the entropy
 305 functional such that in this case, the joint product prior density matrix $\hat{\rho}_1 \otimes \hat{\rho}_2$ updates to the product
 306 posterior $\hat{\rho}_1 \otimes \hat{\rho}_2$ so that inferences about one do not affect inferences about the other.

307 The argument below is considerably simplified if we expand the space of density matrices to
 308 include density matrices that are not necessarily normalized. This does not represent any limitation
 309 because normalization can always be easily achieved as one additional constraint. We consider a few
 310 special cases below:

311 **Case 1:** We receive the extremely constraining information that the posterior distribution for system 1
 is completely specified to be $\hat{\rho}_1$ while we receive no information about system 2 at all. We treat the two
 systems jointly. Maximize the joint entropy $S[\hat{\rho}_{12}, \hat{\rho}_1 \otimes \hat{\rho}_2]$, subject to the following constraints on the
 $\hat{\rho}_{12}$,

$$\text{Tr}_2(\hat{\rho}_{12}) = \hat{\rho}_1. \quad (44)$$

312 Notice all of the N^2 elements in \mathcal{H}_1 of $\hat{\rho}_{12}$ are being constrained. We therefore need a Lagrange
 313 multiplier which spans \mathcal{H}_1 and therefore it is a square matrix $\hat{\lambda}_1$. This is readily seen by observing the
 314 component form expressions of the Lagrange multipliers $(\hat{\lambda}_1)_{ij} = \lambda_{ij}$. Maximizing the entropy with
 315 respect to this \mathcal{H}_2 independent constraint is,

$$0 = \delta \left(S - \sum_{ij} \lambda_{ij} \left(\text{Tr}_2(\hat{\rho}_{1,2}) - \hat{\rho}_1 \right)_{ij} \right), \quad (45)$$

316 but reexpressing this with its transpose $(\hat{\lambda}_1)_{ij} = (\hat{\lambda}_1^T)_{ji}$, gives

$$0 = \delta \left(S - \text{Tr}_1(\hat{\lambda}_1 [\text{Tr}_2(\hat{\rho}_{1,2}) - \hat{\rho}_1]) \right), \quad (46)$$

where we have relabeled $\hat{\lambda}_1^T \rightarrow \hat{\lambda}_1$, for convenience, as the name of the Lagrange multipliers are
 arbitrary. For arbitrary variations of $\hat{\rho}_{12}$, we therefore have,

$$\hat{\lambda}_1 \otimes \hat{1}_2 = \phi(\hat{\rho}_{12}, \hat{\rho}_1 \otimes \hat{\rho}_2). \quad (47)$$

DC2 is implemented by requiring $\hat{\rho}_1 \otimes \hat{\rho}_2 \rightarrow \hat{\rho}_1 \otimes \hat{\rho}_2$, such that the function ϕ is designed to reflect
 subsystem independence in this case; therefore, we have

$$\hat{\lambda}_1 \otimes \hat{1}_2 = \phi(\hat{\rho}_1 \otimes \hat{\rho}_2, \hat{\rho}_1 \otimes \hat{\rho}_2). \quad (48)$$

This equation must hold for all choices of the independent prior $\hat{\rho}_2$ in \mathcal{H}_2 . Suppose we had chosen a
 different prior $\hat{\rho}'_2 = \hat{\rho}_2 + \delta\hat{\rho}_2$. For all $\delta\hat{\rho}_2$ the LHS $\hat{\lambda}_1 \otimes \hat{1}_2$ remains unchanged. This means that any
 dependence that the right hand side might potentially have had on $\hat{\rho}_2$ must cancel out, meaning,

$$\phi(\hat{\rho}_1 \otimes \hat{\rho}_2, \hat{\rho}_1 \otimes \hat{\rho}_2) = f(\hat{\rho}_1, \hat{\rho}_1) \otimes \hat{1}_2. \quad (49)$$

317 Since $\hat{\phi}_2$ is arbitrary, suppose further that we choose a unit prior, $\hat{\phi}_2 = \hat{1}_2$, and note that $\hat{\rho}_1 \otimes \hat{1}_2$ and
 318 $\hat{\rho}_1 \otimes \hat{1}_2$ are block diagonal in \mathcal{H}_2 . Because the LHS is block diagonal in \mathcal{H}_2 ,

$$f(\hat{\rho}_1, \hat{\phi}_1) \otimes \hat{1}_2 = \phi(\hat{\rho}_1 \otimes \hat{1}_2, \hat{\rho}_1 \otimes \hat{1}_2) \quad (50)$$

319 the RHS is block diagonal in \mathcal{H}_2 , and because the function ϕ is understood to be a power series
 320 expansion in its arguments,

$$f(\hat{\rho}_1, \hat{\phi}_1) \otimes \hat{1}_2 = \phi(\hat{\rho}_1 \otimes \hat{1}_2, \hat{\rho}_1 \otimes \hat{1}_2) = \phi(\hat{\rho}_1, \hat{\phi}_1) \otimes \hat{1}_2. \quad (51)$$

This gives,

$$\hat{\lambda}_1 \otimes \hat{1}_2 = \phi(\hat{\rho}_1, \hat{\phi}_1) \otimes \hat{1}_2, \quad (52)$$

321 and therefore the $\hat{1}_2$ factors out and $\hat{\lambda}_1 = \phi(\hat{\rho}_1, \hat{\phi}_1)$. A similar argument exchanging systems 1 and 2
 322 shows $\hat{\lambda}_2 = \phi(\hat{\rho}_2, \hat{\phi}_2)$ in this case.

Case 1 - Conclusion: The analysis leads us to conclude that when the system 2 is not updated the dependence on $\hat{\phi}_2$ also drops out,

$$\phi(\hat{\rho}_1 \otimes \hat{\phi}_2, \hat{\rho}_1 \otimes \hat{\phi}_2) = \phi(\hat{\rho}_1, \hat{\phi}_1) \otimes \hat{1}_2, \quad (53)$$

and similarly,

$$\phi(\hat{\rho}_1 \otimes \hat{\rho}_2, \hat{\rho}_1 \otimes \hat{\rho}_2) = \hat{1}_1 \otimes \phi(\hat{\rho}_2, \hat{\phi}_2). \quad (54)$$

Case 2: Now consider a different special case in which the marginal posterior distributions for systems 1 and 2 are both completely specified to be $\hat{\rho}_1$ and $\hat{\rho}_2$ respectively. Maximize the joint entropy, $S[\hat{\rho}_{12}, \hat{\rho}_1 \otimes \hat{\rho}_2]$, subject to the following constraints on the $\hat{\rho}_{12}$,

$$\text{Tr}_2(\hat{\rho}_{12}) = \hat{\rho}_1 \quad \text{and} \quad \text{Tr}_1(\hat{\rho}_{12}) = \hat{\rho}_2. \quad (55)$$

323 Here each expectation value constraints the entire space \mathcal{H}_i , where $\hat{\rho}_i$ lives. The Lagrange multipliers
 324 must span their respective spaces, so we implement the constraint with the Lagrange multiplier
 325 operator $\hat{\mu}_i$, then,

$$0 = \delta \left(S - \text{Tr}_1(\hat{\mu}_1[\text{Tr}_2(\hat{\rho}_{12}) - \hat{\rho}_1]) - \text{Tr}_2(\hat{\mu}_2[\text{Tr}_1(\hat{\rho}_{12}) - \hat{\rho}_2]) \right). \quad (56)$$

For arbitrary variations of $\hat{\rho}_{12}$, we have,

$$\hat{\mu}_1 \otimes \hat{1}_2 + \hat{1}_1 \otimes \hat{\mu}_2 = \phi(\hat{\rho}_{12}, \hat{\rho}_1 \otimes \hat{\rho}_2). \quad (57)$$

By design, DC2 is implemented by requiring $\hat{\phi}_1 \otimes \hat{\phi}_2 \rightarrow \hat{\rho}_1 \otimes \hat{\rho}_2$ in this case; therefore, we have

$$\hat{\mu}_1 \otimes \hat{1}_2 + \hat{1}_1 \otimes \hat{\mu}_2 = \phi(\hat{\rho}_1 \otimes \hat{\rho}_2, \hat{\rho}_1 \otimes \hat{\rho}_2). \quad (58)$$

Write (58) as,

$$\hat{\mu}_1 \otimes \hat{1}_2 = \phi(\hat{\rho}_1 \otimes \hat{\rho}_2, \hat{\rho}_1 \otimes \hat{\rho}_2) - \hat{1}_1 \otimes \hat{\mu}_2. \quad (59)$$

The left hand side is independent of changes in of $\hat{\rho}_2$ and $\hat{\phi}_2$ in \mathcal{H}_2 as $\hat{\mu}_2$ "pushes out" this dependence from ϕ . Any dependence that the RHS might potentially have had on $\hat{\rho}_2, \hat{\phi}_2$ must cancel out, leaving $\hat{\mu}_1$ unchanged. Consequently,

$$\phi(\hat{\rho}_1 \otimes \hat{\rho}_2, \hat{\rho}_1 \otimes \hat{\rho}_2) - \hat{1}_1 \otimes \hat{\mu}_2 = g(\hat{\rho}_1, \hat{\phi}_1) \otimes \hat{1}_2. \quad (60)$$

Because $g(\hat{\rho}_1, \hat{\phi}_1)$ is independent of arbitrary variations of $\hat{\rho}_2$ and $\hat{\phi}_2$ on the LHS above – it is satisfied equally well for all choices. The form of $g(\hat{\rho}_1, \hat{\phi}_1)$ reduces to the form of $f(\hat{\rho}_1, \hat{\phi}_1)$ from Case 1 when $\hat{\rho}_2 = \hat{\phi}_2 = \hat{1}_2$ and similarly DC1' gives $\hat{\mu}_2 = 0$. Therefore, the Lagrange multiplier is

$$\hat{\mu}_1 \otimes \hat{1}_2 = \phi(\hat{\rho}_1, \hat{\phi}_1) \otimes \hat{1}_2. \quad (61)$$

A similar analysis can be carried out for $\hat{\mu}_2$ leading to

$$\hat{1}_1 \otimes \hat{\mu}_2 = \hat{1}_1 \otimes \phi(\hat{\rho}_2, \hat{\phi}_2). \quad (62)$$

³²⁶ **Case 2 - Conclusion:** Substituting back into (58) gives us a functional equation for ϕ ,

$$\phi(\hat{\rho}_1 \otimes \hat{\rho}_2, \hat{\phi}_1 \otimes \hat{\phi}_2) = \phi(\hat{\rho}_1, \hat{\phi}_1) \otimes \hat{1}_2 + \hat{1}_1 \otimes \phi(\hat{\rho}_2, \hat{\phi}_2), \quad (63)$$

³²⁷ which is,

$$\phi(\hat{\rho}_1 \otimes \hat{\rho}_2, \hat{\phi}_1 \otimes \hat{\phi}_2) = \phi(\hat{\rho}_1 \otimes \hat{1}_2, \hat{\phi}_1 \otimes \hat{1}_2) + \phi(\hat{1}_1 \otimes \hat{\rho}_2, \hat{1}_1 \otimes \hat{\phi}_2). \quad (64)$$

³²⁸ The general solution to this matrix valued functional equation is derived in the Appendix 5.5, and is,

$$\phi(\hat{\rho}, \hat{\phi}) = \tilde{A} \ln(\hat{\rho}) + \tilde{B} \ln(\hat{\phi}), \quad (65)$$

³²⁹ where tilde \tilde{A} is a “super-operator” having constant coefficients and twice the number of indicies as $\hat{\rho}$
³³⁰ and $\hat{\phi}$ as discussed in the Appendix (i.e. $(\tilde{A} \ln(\hat{\rho}))_{ij} = \sum_{k\ell} A_{ijkl} (\log(\hat{\rho}))_{k\ell}$ and similarly for $\tilde{B} \ln(\hat{\phi})$).
³³¹ DC1' imposes,

$$\phi(\hat{\rho}, \hat{\phi}) = \tilde{A} \ln(\hat{\rho}) + \tilde{B} \ln(\hat{\phi}) = \hat{0}, \quad (66)$$

³³² which is satisfied in general when $\tilde{A} = -\tilde{B}$, and now,

$$\phi(\hat{\rho}, \hat{\phi}) = \tilde{A} \left(\ln(\hat{\rho}) - \ln(\hat{\phi}) \right). \quad (67)$$

³³³ We may fix the constant \tilde{A} by substituting our solution into the RHS of equation (63) which is equal to
³³⁴ the RHS of equation (64),

$$\left(\tilde{A}_1 \left(\ln(\hat{\rho}_1) - \ln(\hat{\phi}_1) \right) \right) \otimes \hat{1}_2 + \hat{1}_1 \otimes \left(\tilde{A}_2 \left(\ln(\hat{\rho}_2) - \ln(\hat{\phi}_2) \right) \right)$$

³³⁵

$$= \tilde{A}_{12} \left(\ln(\hat{\rho}_1 \otimes \hat{1}_2) - \ln(\hat{\phi}_1 \otimes \hat{1}_2) \right) + \tilde{A}_{12} \left(\ln(\hat{1}_1 \otimes \hat{\rho}_2) - \ln(\hat{1}_1 \otimes \hat{\phi}_2) \right), \quad (68)$$

³³⁶ where \tilde{A}_{12} acts on the joint space of 1 and 2 and \tilde{A}_1, \tilde{A}_2 acts on single subspaces 1 or 2 respectively.
³³⁷ Using the log tensor product identity, $\ln(\hat{\rho}_1 \otimes \hat{1}_2) = \ln(\hat{\rho}_1) \otimes \hat{1}_2$, in the RHS of equation (68) gives,

$$= \tilde{A}_{12} \left(\ln(\hat{\rho}_1) \otimes \hat{1}_2 - \ln(\hat{\phi}_1) \otimes \hat{1}_2 \right) + \tilde{A}_{12} \left(\hat{1}_1 \otimes \ln(\hat{\rho}_2) - \hat{1}_1 \otimes \ln(\hat{\phi}_2) \right). \quad (69)$$

³³⁸ Note that arbitrarily letting $\hat{\rho}_2 = \hat{\phi}_2$ gives,

$$\left(\tilde{A}_1 \left(\ln(\hat{\rho}_1) - \ln(\hat{\phi}_1) \right) \right) \otimes \hat{1}_2 = \tilde{A}_{12} \left(\ln(\hat{\rho}_1) \otimes \hat{1}_2 - \ln(\hat{\phi}_1) \otimes \hat{1}_2 \right). \quad (70)$$

³³⁹ or arbitrarily letting $\hat{\rho}_1 = \hat{\phi}_1$ gives,

$$\hat{1}_1 \otimes \left(\tilde{A}_2 \left(\ln(\hat{\rho}_2) - \ln(\hat{\phi}_2) \right) \right) = \tilde{A}_{12} \left(\hat{1}_1 \otimes \ln(\hat{\rho}_2) - \hat{1}_1 \otimes \ln(\hat{\phi}_2) \right). \quad (71)$$

³⁴⁰ As \tilde{A}_{12} , \tilde{A}_1 , and \tilde{A}_2 are constant tensors, inspecting the above equalities determines the form of
³⁴¹ the tensor to be $\tilde{A} = A \tilde{1}$ where A is a scalar constant and $\tilde{1}$ is the super-operator identity over the
³⁴² appropriate (joint) Hilbert space.

³⁴³ Because our goal is to maximize the entropy function, we let the arbitrary constant $A = -|A|$ and
³⁴⁴ distribute $\tilde{1}$ identically, which gives the final functional form,

$$\phi(\hat{\rho}, \hat{\phi}) = -|A| \left(\ln(\hat{\rho}) - \ln(\hat{\phi}) \right). \quad (72)$$

³⁴⁵ "Integrating" ϕ , gives a general form for the quantum relative entropy,

$$S(\hat{\rho}, \hat{\phi}) = -|A| \text{Tr}(\hat{\rho} \log \hat{\rho} - \hat{\rho} \log \hat{\phi} - \hat{\rho}) + C[\hat{\phi}] = -|A| S_U(\hat{\rho}, \hat{\phi}) + |A| \text{Tr}(\hat{\rho}) + C[\hat{\phi}], \quad (73)$$

³⁴⁶ where $S_U(\hat{\rho}, \hat{\phi})$ is Umegaki's form of the relative entropy, the extra $|A| \text{Tr}(\hat{\rho})$ from integration is an
³⁴⁷ artifact present for the preservation of DC1', and $C[\hat{\phi}]$ is a constant in the sense that it drops out under
³⁴⁸ arbitrary variations of $\hat{\rho}$. This entropy leads to the same inferences as Umegaki's form of the entropy
³⁴⁹ with added bonus that $\hat{\rho} = \hat{\phi}$ in the absence of constraints or changes in information – rather than
³⁵⁰ $\hat{\rho} = e^{-1} \hat{\phi}$ which would be given by maximizing Umegaki's form of the entropy. In this sense the extra
³⁵¹ $|A| \text{Tr}(\hat{\rho})$ only improves the inference process as it more readily adheres to the PMU though DC1';
³⁵² however now because $S_U \geq 0$, we have $S(\hat{\rho}, \hat{\phi}) \leq \text{Tr}(\hat{\rho}) + C[\hat{\phi}]$, which provides little nuisance. In the
³⁵³ spirit of this derivation we will keep the $\text{Tr}(\hat{\rho})$ term there, but for all practical purposes of inference, as
³⁵⁴ long as there is a normalization constraint, it plays no role, and we find (letting $|A| = 1$ and $C[\hat{\phi}] = 0$),

$$S(\hat{\rho}, \hat{\phi}) \rightarrow S^*(\hat{\rho}, \hat{\phi}) = -S_U(\hat{\rho}, \hat{\phi}) = -\text{Tr}(\hat{\rho} \log \hat{\rho} - \hat{\rho} \log \hat{\phi}), \quad (74)$$

³⁵⁵ Umegaki's form of the relative entropy. $S^*(\hat{\rho}, \hat{\phi})$ is an equally valid entropy because, given
³⁵⁶ normalization is applied, the same selected posterior $\hat{\rho}$ maximizes both $S(\hat{\rho}, \hat{\phi})$ and $S^*(\hat{\rho}, \hat{\phi})$.

³⁵⁷ 3.2. Remarks

³⁵⁸ Due to the universality and the equal application of the PMU by using the same design criteria
³⁵⁹ for both the standard and quantum case, the quantum relative entropy reduces to the standard relative
³⁶⁰ entropy when $[\hat{\rho}, \hat{\phi}] = 0$ or when the experiment being preformed $\hat{\rho} \rightarrow \rho(a) = \text{Tr}(\hat{\rho}|a\rangle\langle a|)$ is known.
³⁶¹ The quantum relative entropy we derive has the correct asymptotic form of the standard relative
³⁶² entropy in the sense of [8–10]. Further connections will be illustrated in a follow up article that is
³⁶³ concerned with direct applications of the quantum relative entropy. Because two entropies are derived
³⁶⁴ in parallel, we expect the well known inferential results and consequences of the relative entropy to
³⁶⁵ have a quantum relative entropy representation.

³⁶⁶ Maximizing the quantum relative entropy with respect to some constraints $\langle \hat{A}_i \rangle$, where $\{\hat{A}_i\}$ are
³⁶⁷ a set of arbitrary Hermitian operators, and normalization $\langle \hat{1} \rangle = 1$, gives the following general solution
³⁶⁸ for the posterior density matrix:

$$\hat{\rho} = \exp \left(\alpha_0 \hat{1} + \sum_i \alpha_i \hat{A}_i + \ln(\hat{\phi}) \right) = \frac{1}{Z} \exp \left(\sum_i \alpha_i \hat{A}_i + \ln(\hat{\phi}) \right) \equiv \frac{1}{Z} \exp \left(\hat{C} \right), \quad (75)$$

³⁶⁹ where α_i are the Lagrange multipliers of the respective constraints and normalization may be factored
³⁷⁰ out of the exponential in general because the identity commutes universally. If $\hat{\phi} \propto \hat{1}$, it is well known
³⁷¹ the analysis arrives at the same expression for $\hat{\rho}$ after normalization as it would if the von Neumann
³⁷² entropy were used, and thus one can find expressions for thermalized quantum states $\hat{\rho} = \frac{1}{Z} e^{-\beta \hat{H}}$. The

373 remaining problem is to solve for the N Lagrange multipliers using their N associated expectation
 374 value constraints. In principle their solution is found by computing Z and using standard methods
 375 from Statistical Mechanics,

$$\langle \hat{A}_i \rangle = -\frac{\partial}{\partial \alpha_i} \ln(Z), \quad (76)$$

376 and inverting to find $\alpha_i = \alpha_i(\langle \hat{A}_i \rangle)$, which has a unique solution due to the joint concavity (convexity
 377 depending on the sign convention) of the quantum relative entropy [8,9] when the constraints are
 378 linear in $\hat{\rho}$. The simple proof that (76) is monotonic in α , and therefore invertible, is that is that its
 379 derivative $\frac{\partial}{\partial \alpha} \langle \hat{A}_i \rangle = \langle \hat{A}_i^2 \rangle - \langle \hat{A}_i \rangle^2 \geq 0$. Between the Zassenhaus formula

$$e^{t(\hat{A}+\hat{B})} = e^{t\hat{A}} e^{t\hat{B}} e^{-\frac{t^2}{2} [\hat{A}, \hat{B}]} e^{\frac{t^3}{6} (2[\hat{B}, [\hat{A}, \hat{B}]] + [\hat{A}, [\hat{A}, \hat{B}]])} \dots, \quad (77)$$

380 and Horn's inequality, the solutions to (76) lack a certain calculational elegance because it is difficult to
 381 express the eigenvalues of $\hat{C} = \log(\hat{\rho}) + \sum \alpha_i \hat{A}_i$ (in the exponential) in simple terms of the eigenvalues
 382 of the \hat{A}_i 's and $\hat{\rho}$, in general, when the matrices do not commute. The solution requires solving the
 383 eigenvalue problem for \hat{C} , such the the exponential of \hat{C} may be taken and evaluated in terms of the
 384 eigenvalues of the $\alpha_i \hat{A}_i$'s and the prior density matrix $\hat{\rho}$. A pedagogical exercise is, starting with a
 385 prior which is a mixture of spin-z up and down $\hat{\rho} = a|+\rangle\langle+| + b|-\rangle\langle-|$ ($a, b \neq 0$) and maximize
 386 the quantum relative entropy with respect to the expectation of a general Hermitian operator. This
 387 example is given in the Appendix 5.6.

388 4. Conclusions:

389 This approach emphasizes the notion that entropy is a tool for performing inference and
 390 downplays counter-notional issues which arise if one interprets entropy as a measure of disorder,
 391 a measure of distinguishability, or an amount of missing information [7]. Because the same design
 392 criteria, guided by the PMU, are applied equally well to the design of a relative and quantum relative
 393 entropy, we find that both the relative and quantum relative entropy are designed for the purpose of
 394 inference. Because the quantum relative entropy is the function which fits the requirements of a tool
 395 designed for inference, we now know what it is and how to use it – formulating an inferential quantum
 396 maximum entropy method. A follow up article is concerned with a few interesting applications of the
 397 quantum maximum entropy method, and in particular it derives the Quantum Bayes Rule.

398 Acknowledgments:**399**

400 I must give ample acknowledgment to Ariel Caticha who suggested the problem of justifying
401 the form of the quantum relative entropy as a criterion for ranking of density matrices. He cleared up
402 several difficulties by suggesting that design constraints be applied to the variational derivative of the
403 entropy rather than the entropy itself. As well, he provided substantial improvements to the method
404 for imposing DC2 that lead to the functional equations for the variational derivatives ($\phi_{12} = \phi_1 + \phi_2$) –
405 with more rigor than in earlier versions of this article. His time and guidance are all greatly appreciated
406 – Thanks Ariel.

407 References

- 408** 1. Shore, J. E.; Johnson, R. W.; Axiomatic derivation of the Principle of Maximum Entropy and the Principle of
409 Minimum Cross-Entropy. *IEEE Trans. Inf. Theory* **1980** *IT-26*, 26-37.
- 410** 2. Shore, J. E.; Johnson, R. W. Properties of Cross-Entropy Minimization. *IEEE Trans. Inf. Theory* **1981** *IT-27*,
411 472-482.
- 412** 3. Csiszár, I. Why least squares and maximum entropy: an axiomatic approach to inference for linear inverse
413 problems. *Ann. Stat.* **1991**, *19*, 2032.
- 414** 4. Skilling, J. The Axioms of Maximum Entropy. *Maximum- Entropy and Bayesian Methods in Science and Engineering*,
415 Dordrecht, Holland, 1988; Erickson, G. J.; Smith, C. R.; Kluwer Academic Publishers: 1988.
- 416** 5. Skilling, J. Classic Maximum Entropy. *Maximum- Entropy and Bayesian Methods in Science and Engineering*,
417 Dordrecht, Holland, 1988; Skilling, J.; Kluwer Academic Publishers: 1988.
- 418** 6. Skilling, J. Quantified Maximum Entropy. *Maximum En- tropy and Bayesian Methods in Science and Engineering*,
419 Dordrecht, Holland, 1988; Fougére, P. F.; Kluwer Academic Publishers: 1990.
- 420** 7. Caticha, A. *Entropic Inference and the Foundations of Physics* (monograph commissioned
421 by the 11th Brazilian Meeting on Bayesian Statistics – EBEB-2012); Available online:
422 <http://www.albany.edu/physics/ACaticha-EIFP-book.pdf>.
- 423** 8. Hiai, F.; Petz, D. The Proper Formula for Relative Entropy and its Asymptotics in Quantum Probability.
424 *Commun. Math. Phys.* **1991**, *143*, 99-114.
- 425** 9. Petz, D. Characterization of the Relative Entropy of States of Matrix Algebras. *Acta Math. Hung.* **1992**, *59*,
426 449-455.
- 427** 10. Ohya, M.; Petz, D. *Quantum Entropy and Its Use*; Springer-Verlag: New York NY, USA, 1993; 0-387-54881-5.
- 428** 11. Wilming, H.; Gallego, R.; Eisert, J. Axiomatic Characterization of the Quantum Relative Entropy and Free
429 Energy. *Entropy* **2017**, *19*, 241, 10.3390/e19060241.
- 430** 12. Jaynes, E.T. Information Theory and Statistical Mechanics. *Phys. Rev.* **1957**, *106*, 620-630.
- 431** 13. Jaynes, E.T. *Probability Theory: The Logic of Science*; Cambridge University Press, Cambridge, UK, 2003.
- 432** 14. Jaynes, E.T. Information Theory and Statistical Mechanics II. *Phys. Rev.* **1957**, *108*, 171-190.
- 433** 15. Balian, R.; Vénéroni, M. Incomplete descriptions, relevant information, and entropy production in collision
434 processes. *Ann. Phys.* **1987**, *174*, 229-224, 10.1016/0003-4916(87)90085-6.
- 435** 16. Balian, R.; Balazs, N. L. Equiprobability, inference and entropy in quantum theory. *Ann. Phys.* **1987**, *179*, *174*,
436 97-144, 10.1016/S0003-4916(87)80006-4.
- 437** 17. Balian, R. Justification of the Maximum Entropy Criterion in Quantum Mechanics. *Maximum Entropy and
438 Bayesian Methods*, 1989, 123-129; Skilling, J.; Kluwer Academic Publishers, 1989., 10.1007/978-94-015-7860-8_9.
- 439** 18. Balian, R., On the principles of quantum mechanics. *Amer. J. Phys.* **1989**, *57*, 1019-1027.
- 440** 19. Balian, R. Gain of information in a quantum measurement. *Eur. J. Phys* **1989**, *10*, 208-213
- 441** 20. Balian, R. Incomplete descriptions and relevant entropies. *Am. J. Phys.* **1999**, *67*, 1078-1090,10.1119/1.19086.
- 442** 21. Blankenbecler, R.; Partovi, H. Uncertainty, Entropy, and the Statistical Mechanics of Microscopic Systems.
443 *Phys. Rev. Lett.* **1985**, *54*, 373-376.
- 444** 22. Blankenbecler, R.; Partovi, H. Quantum Density Matrix and Entropic Uncertainty. *The Fifth Workshop on
445 Maximum Entropy and Bayesian Methods in Applied Statistics*, Laramie WY, USA, August, 1985.
- 446** 23. Korotkov, A. Continuous quantum measurement of a double dot. *Phys. Rev. B* **1999**, *60*, 5737-5742.
- 447** 24. Korotkov, A. Selective quantum evolution of a qubit state due to continuous measurement. *Phys. Rev. B* **2000**,
448 *63*, 115403, 10.1103/PhysRevB.63.115403.

- 449 25. Jordan, A.; Korotkov, A. Qubit feedback and control with kicked quantum nondemolition measurements: A
450 quantum Bayesian analysis. *Phys. Rev. B* **2006**, *74*, 085307.
- 451 26. Hellmann, F.; Kamiński, W.; Kostecki, P. Quantum collapse rules from the maximum relative entropy principle.
452 *New J. Phys.* **2016**, *18*, 013022.
- 453 27. Giffin, A.; Caticha, A. Updating Probabilities. Presented at MaxEnt (2006). *MaxEnt 2006, the 26th International*
- 454 *Workshop on Bayesian Inference and Maximum Entropy Methods*, Paris, France.
- 455 28. Caticha, A. Toward an Informational Pragmatic Realism, *Minds and Machines* **2014**, *24*, 37–70.
- 456 29. Umegaki, H. Conditional expectation in an operator algebra, IV (entropy and information). *Kodai Math. Sem.*
457 *Rep* **1962**, *14*, 59–85.
- 458 30. Uhlmann, A. Relative entropy and the Wigner-Yanase-Dyson-Lieb concavity in an interpolation theory.
459 *Commun. Math. Phys.* **1997**, *54*, 21–32.
- 460 31. Schumacher, B.; Westmoreland, M. Relative entropy in quantum information theory. *AMS special session on*
461 *Quantum Information and Computation*, January, 2000.
- 462 32. von Neumann, J. *Mathematische Grundlagen der Quantenmechanik*; Springer-Verlag: Berlin, Germany, 1932.
463 [English translation: *Mathematical Foundations of Quantum Mechanics*; Princeton University Press, Princeton NY,
464 USA, 1983.]
- 465 33. Aczél, J. *Lectures on Functional Equations and their Applications*, Vol. 19; Academic Press Inc.: 111 Fifth Ave. New
466 York NY 10003, USA, 1966; pp. 31–44, 141–145, 213–217, 301–302, 347–349.

467 **5. Appendix:**

468 The Appendix loosely follows the relevant sections in [33], and then uses the methods reviewed to
469 solve the relevant functional equations for ϕ . The last section is an example of the quantum maximum
470 entropy method for spin.

471 **5.1. Simple functional equations**

472 From [33] pages 31–44.

473 Thm 1:

474 If Cauchy's functional equation

$$f(x+y) = f(x) + f(y), \quad (78)$$

475 is satisfied for all real x, y , and if the function $f(x)$ is (a) continuous at a point, (b) nonnegative for small positive
476 x 's, or (c) bounded in an interval, then,

$$f(x) = cx \quad (79)$$

477 is the solution to (78) for all real x . If (78) is assumed only over all positive x, y , then under the same conditions
478 (79) holds for all positive x .

479 Proof 1:

480 The most natural assumption for our purposes is that $f(x)$ is continuous at a point (which later
481 extends to continuity all points as given by Darboux). Cauchy solved the functional equation by
482 induction. In particular equation (78) implies,

$$f\left(\sum_i x_i\right) = \sum_i f(x_i), \quad (80)$$

483 and if we let each $x_i = x$ as a special case to determine f , we find

$$f(nx) = nf(x). \quad (81)$$

484 We may let $nx = mt$ such that

$$f(x) = f\left(\frac{m}{n}t\right) = \frac{m}{n}f(t). \quad (82)$$

485 Letting $\lim_{t \rightarrow 1} f(t) = f(1) = c$, gives

$$f\left(\frac{m}{n}\right) = \frac{m}{n}f(1) = \frac{m}{n}c, \quad (83)$$

486 and because for $t = 1$, $x = \frac{m}{n}$ above, we have

$$f(x) = cx, \quad (84)$$

487 which is the general solution of the linear functional equation. In principle c can be complex. The
488 importance of Cauchy's solution is that can be used to give general solutions to the following Cauchy
489 equations:

$$f(x+y) = f(x)f(y), \quad (85)$$

$$f(xy) = f(x) + f(y), \quad (86)$$

$$f(xy) = f(x)f(y), \quad (87)$$

490 by performing consistent substitution until they are the same form as (78) as given by Cauchy. We will
491 briefly discuss the first two.

492 Thm 2:

493 The general solution of $f(x+y) = f(x)f(y)$ is $f(x) = e^{cx}$ for all real or for all positive x, y that are
494 continuous at one point and, in addition to the exponential solution, the solution $f(0) = 1$ and $f(x) = 0$ for
495 ($x > 0$) are in these classes of functions.

496 The first functional $f(x+y) = f(x)f(y)$ is solved by first noting that it is strictly positive for real
497 $x, y, f(x)$, which can be shown by considering $x = y$,

$$f(2x) = f(x)^2 > 0. \quad (88)$$

498 If there exists $f(x_0) = 0$, then it follows that $f(x) = f((x - x_0) + x_0) = 0$, a trivial solution, hence why
499 the possibility of being equal to zero is excluded above. Given $f(x)$ is nowhere zero, we are justified in
500 taking the natural logarithm $\ln(x)$, due to its positivity $f(x) > 0$. This gives,

$$\ln(f(x+y)) = \ln(f(x)) + \ln(f(y)), \quad (89)$$

501 and letting $g(x) = \ln(f(x))$ gives,

$$g(x+y) = g(x) + g(y), \quad (90)$$

502 which is Cauchy's linear equation, and thus has the solution $g(x) = cx$. Because $g(x) = \ln(f(x))$, one
503 finds in general that $f(x) = e^{cx}$.

504 Thm 3:

505 If the functional equation $f(xy) = f(x) + f(y)$ is valid for all positive x, y then its general solution is
506 $f(x) = c \ln(x)$ given it is continuous at a point. If $x = 0$ (or $y = 0$) are valid then the general solution is
507 $f(x) = 0$. If all real x, y are valid except 0 then the general solution is $f(x) = c \ln(|x|)$.

508 In particular we are interested in the functional equation $f(xy) = f(x) + f(y)$ when x, y are
 509 positive. In this case we can again follow Cauchy and substitute $x = e^u$ and $y = e^v$ to get,

$$f(e^u e^v) = f(e^u) + f(e^v), \quad (91)$$

510 and letting $g(u) = f(e^u)$ gives $g(u + v) = g(u) + g(v)$. Again, the solution is $g(u) = cu$ and
 511 therefore the general solution is $f(x) = c \ln(x)$ when we substitute for u . If x could equal 0 then
 512 $f(0) = f(x) + f(0)$, which has the trivial solution $f(x) = 0$. The general solution for $x \neq 0, y \neq 0$ and
 513 x, y positive is therefore $f(x) = c \ln(x)$.

514 *5.2. Functional equations with multiple arguments*

515 From [33] pages 213-217. Consider the functional equation,

$$F(x_1 + y_1, x_2 + y_2, \dots, x_n + y_n) = F(x_1, x_2, \dots, x_n) + F(y_1, y_2, \dots, y_n), \quad (92)$$

516 which is a generalization of Cauchy's linear functional equation (78) to several arguments. Letting
 517 $x_2 = x_3 = \dots = x_n = y_2 = y_3 = \dots = y_n = 0$ gives

$$F(x_1 + y_1, 0, \dots, 0) = F(x_1, 0, \dots, 0) + F(y_1, 0, \dots, 0), \quad (93)$$

518 which is the Cauchy linear functional equation having solution $F(x_1, 0, \dots, 0) = c_1 x_1$ where $F(x_1, 0, \dots, 0)$
 519 is assumed to be continuous or at least measurable majorant. Similarly,

$$F(0, \dots, 0, x_k, 0, \dots, 0) = c_k x_k, \quad (94)$$

520 and if you consider

$$F(x_1 + 0, 0 + y_2, 0, \dots, 0) = F(x_1, 0, \dots, 0) + F(0, y_2, 0, \dots, 0) = c_1 x_1 + c_2 y_2, \quad (95)$$

521 and as y_2 is arbitrary we could have let $y_2 = x_2$ such that in general

$$F(x_1, x_2, \dots, x_n) = \sum c_i x_i, \quad (96)$$

522 as a general solution.

523 *5.3. Relative entropy:*

524 We are interested in the following functional equation,

$$\phi(\rho_1 \rho_2, \varphi_1 \varphi_2) = \phi(\rho_1, \varphi_1) + \phi(\rho_2, \varphi_2). \quad (97)$$

525 This is an equation of the form,

$$F(x_1 y_1, x_2 y_2) = F(x_1, x_2) + F(y_1, y_2), \quad (98)$$

526 where $x_1 = \rho(x_1)$, $y_1 = \rho(x_2)$, $x_2 = \varphi(x_1)$, and $y_2 = \varphi(x_2)$. First assume all ρ and φ are greater than
 527 zero. Then, substitute: $x_i = e^{x'_i}$ and $y_i = e^{y'_i}$ and let $F'(x'_1, x'_2) = F(e^{x'_1}, e^{x'_2})$ and so on such that

$$F'(x'_1 + y'_1, x'_2 + y'_2) = F'(x'_1, x'_2) + F'(y'_1, y'_2), \quad (99)$$

528 which is of the form of (92). The general solution for F is therefore

$$F'(x'_1 + y'_1, x'_2 + y'_2) = a_1(x'_1 + y'_1) + a_2(x'_2 + y'_2) = a_1 \ln(x_1 y_1) + a_2 \ln(x_2 y_2) = F(x_1 y_1, x_2 y_2) \quad (100)$$

529 which means the general solution for ϕ is,

$$\phi(\rho_1, \varphi_1) = a_1 \ln(\rho(x_1)) + a_2 \ln(\varphi(x_1)) \quad (101)$$

530 In such a case when $\varphi(x_0) = 0$ for some value $x_0 \in \mathcal{X}$ we may let $\varphi(x_0) = \epsilon$ where ϵ is as close to zero
531 as we could possibly want – the trivial general solution $\phi = 0$ is saturated by the special case when
532 $\rho = \varphi$ from DC1'. Here we return to the text.

533 *5.4. Matrix functional equations*

534 (This derivation is implied in [33] pages 347-349). First consider a Cauchy matrix functional
535 equation,

$$f(\hat{X} + \hat{Y}) = f(\hat{X}) + f(\hat{Y}) \quad (102)$$

536 where \hat{X} and \hat{Y} are $n \times n$ square matrices. Rewriting the matrix functional equation in terms of its
537 components gives,

$$f_{ij}(x_{11} + y_{11}, x_{12} + y_{12}, \dots, x_{nn} + y_{nn}) = f_{ij}(x_{11}, x_{12}, \dots, x_{nn}) + f_{ij}(y_{11}, y_{12}, \dots, y_{nn}) \quad (103)$$

538 is now in the form of (92) and therefore the solution is,

$$f_{ij}(x_{11}, x_{12}, \dots, x_{nn}) = \sum_{\ell,k=0}^n c_{ij\ell k} x_{\ell k} \quad (104)$$

539 for $i, j = 1, \dots, n$. We find it convenient to introduce super indices, $A = (i, j)$ and $B = (\ell, k)$ such that
540 the component equation becomes,

$$f_A = \sum_B c_{AB} x_B. \quad (105)$$

541 resembles the solution for a linear transformation of a vector from [33]. In general we will be discussing
542 matrices $\hat{X} = \hat{X}_1 \otimes \hat{X}_2 \otimes \dots \otimes \hat{X}_N$ which stem out of the tensor products of density matrices. In this
543 situation \hat{X} can be thought of as $2N$ index tensor or a $z \times z$ matrix where $z = \prod_i^N n_i$ is the product of
544 the ranks of the matrices in the tensor product or even \hat{X} is a vector of length z^2 . In such a case we may
545 abuse the super index notation where A and B lump together the appropriate number of indices such
546 that (105) is the form of the solution for the components in general. The matrix form of the general
547 solution is,

$$f(\hat{X}) = \tilde{C} \hat{X}, \quad (106)$$

548 where \tilde{C} is a constant super-operator having components c_{AB} .

549 *5.5. Quantum Relative entropy:*

550 The functional equation is,

$$\phi(\hat{\rho}_1 \otimes \hat{\rho}_2, \hat{\varphi}_1 \otimes \hat{\varphi}_2) = \phi(\hat{\rho}_1 \otimes \hat{1}_2, \hat{\varphi}_1 \otimes \hat{1}_2) + \phi(\hat{1}_1 \otimes \hat{\rho}_2, \hat{1}_1 \otimes \hat{\varphi}_2). \quad (107)$$

551 These density matrices are Hermitian, positive semi-definite, have positive eigenvalues, and are not
552 equal to $\hat{0}$. Because every invertible matrix can be expressed as the exponential of some other matrix,
553 we can substitute $\hat{\rho}_1 = e^{\hat{\rho}'_1}$, and so on for all four density matrices which gives,

$$\phi(e^{\hat{\rho}'_1} \otimes e^{\hat{\rho}'_2}, e^{\hat{\varphi}'_1} \otimes e^{\hat{\varphi}'_2}) = \phi(e^{\hat{\rho}'_1} \otimes \hat{1}_2, e^{\hat{\varphi}'_1} \otimes \hat{1}_2) + \phi(\hat{1}_1 \otimes e^{\hat{\rho}'_2}, \hat{1}_1 \otimes e^{\hat{\varphi}'_2}). \quad (108)$$

554 Now we use the following identities for Hermitian matrices,

$$e^{\hat{\rho}'_1} \otimes e^{\hat{\rho}'_2} = e^{\hat{\rho}'_1 \otimes \hat{1}_2 + \hat{1}_1 \otimes \hat{\rho}'_2} \quad (109)$$

555 and

$$e^{\hat{\rho}'_1} \otimes \hat{1}_2 = e^{\hat{\rho}'_1 \otimes \hat{1}_2}, \quad (110)$$

556 to recast the functional equation as,

$$\phi\left(e^{\hat{\rho}'_1 \otimes \hat{1}_2 + \hat{1}_1 \otimes \hat{\rho}'_2}, e^{\hat{\rho}'_1 \otimes \hat{1}_2 + \hat{1}_1 \otimes \hat{\rho}'_2}\right) = \phi\left(e^{\hat{\rho}'_1 \otimes \hat{1}_2}, e^{\hat{\rho}'_1 \otimes \hat{1}_2}\right) + \phi\left(e^{\hat{1}_1 \otimes \hat{\rho}'_2}, e^{\hat{1}_1 \otimes \hat{\rho}'_2}\right). \quad (111)$$

557 Letting $G(\hat{\rho}'_1 \otimes \hat{1}_2, \hat{\rho}'_1 \otimes \hat{1}_2) = \phi\left(e^{\hat{\rho}'_1 \otimes \hat{1}_2}, e^{\hat{\rho}'_1 \otimes \hat{1}_2}\right)$ gives,

$$G(\hat{\rho}'_1 \otimes \hat{1}_2 + \hat{1}_1 \otimes \hat{\rho}'_2, \hat{\rho}'_1 \otimes \hat{1}_2 + \hat{1}_1 \otimes \hat{\rho}'_2) = G(\hat{\rho}'_1 \otimes \hat{1}_2, \hat{\rho}'_1 \otimes \hat{1}_2) + G(\hat{1}_1 \otimes \hat{\rho}'_2, \hat{1}_1 \otimes \hat{\rho}'_2). \quad (112)$$

558 This functional equation is of the form

$$G(\hat{X}'_1 + \hat{Y}'_1, \hat{X}'_2 + \hat{Y}'_2) = G(\hat{X}'_1, \hat{X}'_2) + G(\hat{Y}'_1, \hat{Y}'_2), \quad (113)$$

559 which has the general solution

$$G(\hat{X}', \hat{Y}') = \tilde{A} \hat{X}' + \tilde{B} \hat{Y}', \quad (114)$$

560 synonymous to (96), and finally in general,

$$\phi(\hat{\rho}, \hat{\phi}) = \tilde{A} \ln(\hat{\rho}) + \tilde{B} \ln(\hat{\phi}). \quad (115)$$

561 where \tilde{A}, \tilde{B} are super-operators having constant coefficients.

562 5.6. Spin Example

563 Consider an arbitrarily mixed prior is (in the spin- z basis for convenience) with $a, b \neq 0$,

$$\hat{\phi} = a|+\rangle\langle+| + b|-\rangle\langle-| \quad (116)$$

564 and a general Hermitian matrix in the spin-1/2 Hilbert space,

$$c_\mu \hat{\sigma}^\mu = c_1 \hat{1} + c_x \hat{\sigma}_x + c_y \hat{\sigma}_y + c_z \hat{\sigma}_z \quad (117)$$

565

$$= (c_1 + c_z)|+\rangle\langle+| + (c_x - ic_y)|+\rangle\langle-| + (c_x + ic_y)|-\rangle\langle+| + (c_1 - c_z)|-\rangle\langle-|, \quad (118)$$

566 having a known expectation value,

$$\text{Tr}(\hat{\rho} c_\mu \hat{\sigma}^\mu) = c. \quad (119)$$

567 Maximizing the entropy with respect to this general expectation value and normalization is:

$$0 = \left(\delta S - \lambda [\text{Tr}(\hat{\rho}) - 1] - \alpha (\text{Tr}(\hat{\rho} c_\mu \hat{\sigma}^\mu) - c) \right), \quad (120)$$

568 which after varying gives,

$$\hat{\rho} = \frac{1}{Z} \exp(\alpha c_\mu \hat{\sigma}^\mu + \log(\hat{\varphi})). \quad (121)$$

569 Letting

$$\hat{C} = \alpha c_\mu \hat{\sigma}^\mu + \log(\hat{\varphi}) \quad (122)$$

570 gives

$$\begin{aligned} \hat{\rho} &= \frac{1}{Z} e^{\hat{C}} = U e^{U^{-1} \hat{C} U} U^{-1} = \frac{1}{Z} U e^{\hat{\lambda}} U^{-1} \\ &= \frac{e^{\lambda_+}}{Z} U |\lambda_+ \rangle \langle \lambda_+ | U^{-1} + \frac{e^{\lambda_-}}{Z} U |\lambda_- \rangle \langle \lambda_- | U^{-1}, \end{aligned} \quad (123)$$

571 where $\hat{\lambda}$ is the diagonalized matrix of \hat{C} having the real eigenvalues. They are,

$$\lambda_\pm = \lambda \pm \delta\lambda, \quad (124)$$

572 due to the quadratic formula, explicitly:

$$\lambda = \alpha c_1 + \frac{1}{2} \log(ab), \quad (125)$$

573 and

$$\delta\lambda = \frac{1}{2} \sqrt{\left(2\alpha c_z + \log\left(\frac{a}{b}\right)\right)^2 + 4\alpha^2(c_x^2 + c_y^2)}. \quad (126)$$

574 Because λ_\pm and a, b, c_1, c_x, c_y, c_z are real, $\delta\lambda \geq 0$. The normalization constraint specifies the Lagrange
575 multiplier Z ,

$$1 = \text{Tr}(\hat{\rho}) = \frac{e^{\lambda_+} + e^{\lambda_-}}{Z}, \quad (127)$$

576 so $Z = e^{\lambda_+} + e^{\lambda_-} = 2e^\lambda \cosh(\delta\lambda)$. The expectation value constraint specifies the Lagrange multiplier
577 α ,

$$c = \text{Tr}(\hat{\rho} c_\mu \sigma^\mu) = \frac{\partial}{\partial \alpha} \log(Z) = c_1 + \tanh(\delta\lambda) \frac{\partial}{\partial \alpha} \delta\lambda, \quad (128)$$

578 which becomes

$$c = c_1 + \frac{\tanh(\delta\lambda)}{2\delta\lambda} \left(2\alpha(c_x^2 + c_y^2 + c_z^2) + c_z \log\left(\frac{a}{b}\right) \right),$$

579 or

$$c = c_1 + \tanh \left(\frac{1}{2} \sqrt{\left(2\alpha c_z + \log\left(\frac{a}{b}\right)\right)^2 + 4\alpha^2(c_x^2 + c_y^2)} \right) \frac{2\alpha(c_x^2 + c_y^2 + c_z^2) + c_z \log\left(\frac{a}{b}\right)}{\sqrt{\left(2\alpha c_z + \log\left(\frac{a}{b}\right)\right)^2 + 4\alpha^2(c_x^2 + c_y^2)}}. \quad (129)$$

580 This equation is monotonic in α and therefore it is uniquely specified by the value of c . Ultimately this
581 is a consequence from the concavity of the entropy. The proof of (129)'s monotonicity is below:

582 Proof:

583 For $\hat{\rho}$ to be Hermitian, \hat{C} is Hermitian and $\delta\lambda = \frac{1}{2}\sqrt{f(\alpha)}$ is real. Further more, because $\delta\lambda$ is real
 584 $f(\alpha) \geq 0$ and thus $\delta\lambda \geq 0$. Because $f(\alpha)$ is quadratic in α and positive, it may be written in vertex
 585 form,

$$f(\alpha) = a(\alpha - h)^2 + k, \quad (130)$$

586 where $a > 0$, $k \geq 0$, and (h, k) are the (x, y) coordinates of the minimum of $f(\alpha)$. Notice that the form
 587 of (129) is,

$$F(\alpha) = \frac{\tanh(\frac{1}{2}\sqrt{f(\alpha)})}{\sqrt{f(\alpha)}} \times \frac{\partial f(\alpha)}{\partial \alpha}. \quad (131)$$

588 Making the change of variables $\alpha' = \alpha - h$ centers the function such that $f(\alpha') = f(-\alpha')$ is symmetric
 589 about $\alpha' = 0$. We can then write,

$$F(\alpha') = \frac{\tanh(\frac{1}{2}\sqrt{f(\alpha')})}{\sqrt{f(\alpha')}} \times 2a\alpha', \quad (132)$$

590 where the derivative has been computed. Because $f(\alpha')$ is a positive, symmetric, and monotonically
 591 increasing on the (symmetric) half-plane (for α' greater than or less than zero), $S(\alpha') \equiv \frac{\tanh(\frac{1}{2}\sqrt{f(\alpha')})}{\sqrt{f(\alpha')}}$ is
 592 also positive and symmetric, but it is unclear whether or not $S(\alpha)$ is also monotonic in the half-plane.
 593 We may restate

$$F(\alpha') = S(\alpha') \times 2a\alpha'. \quad (133)$$

594 We are now in a decent position to preform the derivate test for monotonic functions:

$$\begin{aligned} \frac{\partial}{\partial \alpha'} F(\alpha') &= 2aS(\alpha') + 2a\alpha' \frac{\partial}{\partial \alpha'} S(\alpha') \\ &= 2aS(\alpha') \left(1 - \frac{a\alpha'^2}{a\alpha'^2 + k}\right) + a \frac{a\alpha'^2}{a\alpha'^2 + k} \left(1 - \tanh^2\left(\frac{1}{2}\sqrt{a\alpha'^2 + k}\right)\right) \\ &\geq 2aS(\alpha') \left(1 - \frac{a(\alpha')^2}{a\alpha'^2 + k}\right) \geq 0 \end{aligned} \quad (134)$$

595 because $a, k, S(\alpha')$, and therefore $\frac{a\alpha'^2}{a\alpha'^2 + k}$ are all > 0 . The function of interest $F(\alpha')$ is therefore monotonic
 596 for all α' , and therefore it is monotonic for all α , completing the proof that there exists a unique real
 597 Lagrange multiplier α in (129).

598 Although (129) is monotonic in α it is seemingly a transcendental equation. This can be solved
 599 graphically for the given values c, c_1, c_x, c_y, c_z , i.e. given the Hermitian matrix and its expectation value
 600 are specified. Equation (129) and the eigenvalues take a simpler form when $a = b = \frac{1}{2}$, because in this
 601 instance $\hat{\phi} \propto \hat{1}$ and commutes universally so it may be factored out of the exponential in (121).