

Bayesian generalized probability calculus for density matrices

Manfred K. Warmuth · Dima Kuzmin

Received: 29 December 2008 / Revised: 12 May 2009 / Accepted: 2 June 2009 /
Published online: 23 July 2009
© The Author(s) 2009. This article is published with open access at Springerlink.com

Abstract One of the main concepts in quantum physics is a density matrix, which is a symmetric positive definite matrix of trace one. Finite probability distributions can be seen as a special case when the density matrix is restricted to be diagonal.

We develop a probability calculus based on these more general distributions that includes definitions of joints, conditionals and formulas that relate these, including analogs of the Theorem of Total Probability and various Bayes rules for the calculation of posterior density matrices. The resulting calculus parallels the familiar “conventional” probability calculus and always retains the latter as a special case when all matrices are diagonal. We motivate both the conventional and the generalized Bayes rule with a minimum relative entropy principle, where the Kullback-Leibler version gives the conventional Bayes rule and Umegaki’s quantum relative entropy the new Bayes rule for density matrices.

Whereas the conventional Bayesian methods maintain uncertainty about which model has the highest data likelihood, the generalization maintains uncertainty about which unit direction has the largest variance. Surprisingly the bounds also generalize: as in the conventional setting we upper bound the negative log likelihood of the data by the negative log likelihood of the MAP estimator.

Keywords Generalized probability · Probability calculus · Density matrix · Quantum Bayes rule

Editor: Nicolo Cesa-Bianchi.

Supported by NSF grant IIS 0325363. Some of this work was done while visiting National ICT Australia in Canberra.

M.K. Warmuth (✉) · D. Kuzmin
UC California—Santa Cruz, Santa Cruz, CA 95064, USA
e-mail: manfred@cse.ucsc.edu

D. Kuzmin
e-mail: dimakuzmin@google.com

1 Introduction

The main notion of a “mixture state” used in quantum physics is a density matrix. States are unit vectors \mathbf{u} ($\|\mathbf{u}\|_2 = 1$). For the sake of simplicity we assume in this paper that the underlying vector space is \mathbb{R}^n (for finite n). Each state \mathbf{u} (unit column vector in \mathbb{R}^n) is associated with a *dyad* $\mathbf{u}\mathbf{u}^\top \in \mathbb{R}^{n \times n}$. The dyad $\mathbf{u}\mathbf{u}^\top$ may be seen as a one-dimensional projection matrix which projects any vector onto direction \mathbf{u} . These dyads are the elementary events of a *generalized probability space*. It is useful to keep the corresponding “conventional” probability space in mind, which consists of a finite set of size n . The n points are the elementary events and a probability distribution may be seen as a mixture over the n points, i.e. such a probability distribution is specified by n real numbers that are bigger than zero and add to one. In the generalized case there are infinitely many dyads even if the dimension n is finite.¹

Density matrices generalize finite probability distributions. They can be defined as mixtures of dyads $\mathbf{W} = \sum_i \omega_i \mathbf{w}_i \mathbf{w}_i^\top$ where the mixture coefficients ω_i are non-negative and sum to one. There may be an arbitrary number of components in the mixture. However, any n dimensional density matrix can be decomposed into a mixture of n orthogonal *eigendyads*, one for each eigenvector (see Fig. 1). Mixtures of dyads are always symmetric² and positive definite. A density matrix \mathbf{W} can be depicted as an ellipse which is an affine transformation of the unit ball: $\{\mathbf{W}\mathbf{u} : \|\mathbf{u}\|_2 = 1\}$ (see Fig. 2). A dyad is a degenerate ellipse with a single axis in direction $\pm\mathbf{u}$ that has radius one (Fig. 1). Note that dyads have trace one:

$$\text{tr}(\mathbf{u}\mathbf{u}^\top) = \text{tr}(\mathbf{u}^\top \mathbf{u}) = \|\mathbf{u}\|_2^2 = 1.$$

Therefore, density matrices also have trace one.

A density matrix \mathbf{W} assigns generalized probability $\text{tr}(\mathbf{W}\mathbf{u}\mathbf{u}^\top)$ to each unit vector \mathbf{u} and its associated dyad $\mathbf{u}\mathbf{u}^\top$ (see Fig. 2). This probability is independent of how \mathbf{W} is expressed as a mixture and can be rewritten as $\mathbf{u}^\top \mathbf{W} \mathbf{u}$. Note that if the symmetric positive definite matrix \mathbf{A} is viewed as a covariance matrix of a random cost vector \mathbf{c} , then $\mathbf{u}^\top \mathbf{A} \mathbf{u}$ is the variance of the cost along direction \mathbf{u} , i.e. the variance of $\mathbf{c} \cdot \mathbf{u}$.

If $\boldsymbol{\alpha} = (\alpha_1, \dots, \alpha_n)$ is a probability vector, then the n -dimensional matrix $\text{diag}(\boldsymbol{\alpha})$ with vector $\boldsymbol{\alpha}$ as its diagonal is a density matrix. Note that $\text{diag}(\boldsymbol{\alpha}) = \sum_i \alpha_i \mathbf{e}_i \mathbf{e}_i^\top$, where the \mathbf{e}_i are the standard basis vectors. Thus conventional probability distributions are special density matrices where the eigensystem is restricted to be the identity matrix. In this paper we develop a Bayesian style analysis for the case when the eigensystem is allowed to be arbitrary.

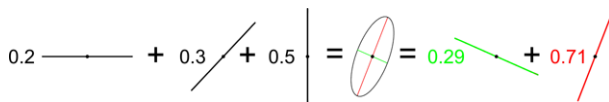


Fig. 1 Two different dyad mixtures that lead to the same density matrix: $0.2 \begin{pmatrix} 1 & 0 \\ 0 & 0 \end{pmatrix} + 0.3 \begin{pmatrix} \sqrt{2}/2 & \sqrt{2}/2 \\ \sqrt{2}/2 & \sqrt{2}/2 \end{pmatrix} + 0.5 \begin{pmatrix} 0 & 0 \\ 0 & 1 \end{pmatrix} = \begin{pmatrix} 0.35 & 0.15 \\ 0.15 & 0.65 \end{pmatrix} = 0.29 \begin{pmatrix} -0.92 & 0.38 \\ 0.38 & -0.92 \end{pmatrix} + 0.71 \begin{pmatrix} 0.38 & 0.92 \\ 0.92 & 0.38 \end{pmatrix}$. Matrices are depicted as ellipses and dyads are degenerate single axis ellipses

¹The machinery for infinite dimensional vector spaces is available. However, in this paper we start with the simplest finite dimensional setting.

²In quantum physics complex numbers are used instead of reals. In that case “symmetric” is replaced by “hermitian” and all our formulas hold for that case as well.

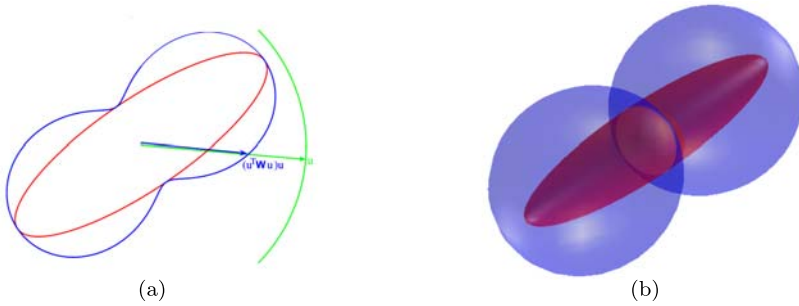


Fig. 2 (Color online) Figure (a) depicts a red ellipse $\{Wu : \|u\|_2 = 1\}$ for some density matrix W . The green curve shows part of the unit ball. The blue figure-eight is a plot of the generalized probabilities in direction u , i.e. $\text{tr}(Wuu^\top)u$. Figure (b) plots a 3-dimensional density matrix (red ellipsoid) and its associated generalized probability surface (in blue)

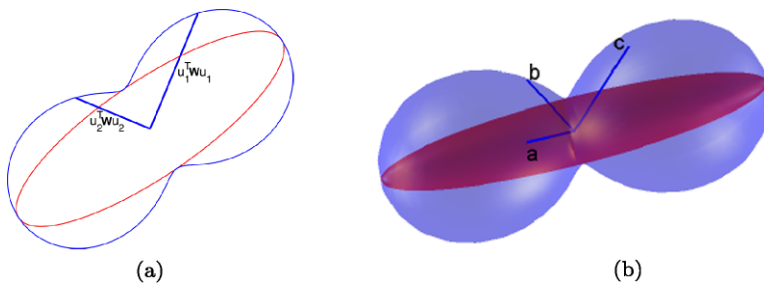


Fig. 3 (Color online) For a set of orthogonal directions u_i and a density matrix W , the sum of generalized probabilities $\text{tr}(Wu_iu_i^\top)$ over the set is one. Figure (a) shows this for 2-dimensional case: the red ellipse is a density matrix W , the blue figure-eight is a plot of the generalized probability $\text{tr}(Wuu^\top)$ around the circle, and for any two orthogonal vectors u_1 and u_2 , $\text{tr}(Wu_1u_1^\top) + \text{tr}(Wu_2u_2^\top) = 1$. Figure (b) shows the three-dimensional case: for any three orthogonal directions u_1 , u_2 and u_3 , the probabilities a , b and c of the three associated dyads sum to one

Perhaps the simplest case to see that something unusual is going on is the uniform density matrix, i.e. $\frac{1}{n}$ times identity I . This density matrix assigns probability $\frac{1}{n}$ to every unit vector, even though there are infinitely many of them. However, note that the sum of generalized probabilities of any set of n orthogonal dyads is $n \frac{1}{n} = 1$. As a matter of fact for any density matrix W and any set of n orthogonal directions u_i , the total generalized probability is one (see Fig. 3)

$$\sum_{i=1}^n \text{tr}(Wu_iu_i^\top) = \text{tr}\left(W \underbrace{\sum_i u_iu_i^\top}_I\right) = \text{tr}(W) = 1. \quad (1.1)$$

This means that while in the conventional case probabilities are additive over the points in the set, in the generalized case probabilities are additive over orthogonal sets of dyads.

In this paper we use density matrices as generalized priors and develop a unifying Bayesian probability calculus for density matrices with rules for translating between joints and conditionals. All formulas retain the conventional case as the special case when the matrices are diagonal. In previous work (Warmuth 2005) we derived a generalized Bayes rule based on the minimum relative entropy principle, but no satisfactory probabilistic interpre-

tation was given for this rule. This Bayes rule fits nicely into our new calculus and we can interpret it using the notion of generalized probability introduced above.

For any fixed orthonormal system \mathbf{u}_i , one can use the dyads $\mathbf{u}_i \mathbf{u}_i^\top$ as elementary events of a conventional probability space. As already discussed, any density matrix can be seen as assigning conventional probabilities to these events that sum to one. Thus if the orthonormal system is fixed, generalized probability space is reduced to conventional probability space over the vectors in the chosen system. Our approach is fundamentally different in that we use density matrices to maintain uncertainty over all orthonormal systems. Our conditional density matrices are part of the probabilistic system specified by a generalized joint probability distribution. In particular, our conditioning method leads to generalizations of the theorem of total probability that involve density matrices.

In Tsuda et al. (2005) various on-line learning updates were generalized from vector parameters to matrix parameters. Following (Kivinen and Warmuth 1997), the updates were derived by minimizing the loss on the current instance plus a divergence to the last parameter. In this paper we use the same method for deriving a Bayes rule for density matrices, which becomes the foundation of our generalized probability calculus. When the parameters are probability vectors over the set of models, then the “conventional” Bayes rule can be derived using the relative entropy as the divergence (e.g. Zellner 1998; Kivinen and Warmuth 1999; Singh et al. 2003). Analogously, we now use the quantum relative entropy, introduced by Umegaki, to derive the generalized Bayes rule.

The new rule uses matrix logarithms and exponentials to avoid the fact that symmetric positive definite matrices are not closed under the matrix product. The rule is strikingly similar to the conventional Bayes rule and retains the latter as a special case when the matrices are diagonal. Various cancellations occur when the conventional Bayes rule is applied iteratively and as we shall see, similar cancellations happen with the new rule (see Sect. 9.2). The conventional Bayes rule may be seen as a soft maximum calculation and the new rule as a soft calculation of the eigenvector with the largest eigenvalue (see Figs. 7 and 8). In Figs. 9 and 10 we plot the projections of posterior onto the eigendirections of the fixed data-likelihood matrix $\mathbf{D}(y|\mathbb{M})$. The projection onto the eigendirection of the largest eigenvalue is a sigmoid like function.

The mathematics applied in this paper are most commonly used in quantum physics. For example, the assignment of generalized probabilities $\text{tr}(\mathbf{W}\mathbf{u}\mathbf{u}^\top)$, can be seen as the outcome of a quantum measurement of a system in mixture state \mathbf{W} being acted upon by a measurement apparatus described by the dyad $\mathbf{u}\mathbf{u}^\top$. It is tempting to call the new rule the “quantum Bayes rule”. However, we currently do not have a quantum physical interpretation of this rule. In particular, the state collapse following a measurement does not explicitly appear in our calculus, also our Bayes rule can not be described as a unitary evolution of the prior state. The term “quantum Bayes rule” also has been claimed before in Schack et al. (2001), where they derive a rule that describes uncertainty information about unobserved quantum measurements of a composite system as a density matrix.

Our work is most closely related to a paper by Cerf and Adami (1999), where, in the context of quantum information theory, a formula was proposed for the conditional density matrix that uses the matrix exponential and matrix logarithm. This special formula appears in our calculus and is now put in a more general context. We hope to transfer many techniques developed in Bayesian Statistics based on the conventional Bayes rule to the context of generalized probabilities.

The use of the quantum relative entropy as a regularizer is reminiscent of the work on quantum state estimation (Olivares and Paris 2007; Bužek et al. 1999). Curiously, the update rules produced in that work do not make use of the matrix logarithm and matrix exponentials

as our new Bayes rule does. Conceivably our calculus can be applied to the state estimation problem.

The paper is organized as follows. Section 2 recalls the relevant matrix algebra facts. Section 3 introduces density matrices and generalized probability distributions and states Gleason's theorem that establishes an equivalence between them. Then, in Sect. 4 we introduce a generalization \odot of the matrix product that is commutative and preserves positive definiteness. This \odot operation is central to our calculus. Section 5 introduces generalized joint distributions. Section 6 discusses marginalizing the joints. Next, in Sect. 7 we give formulas for conditional density matrices. Section 8 presents generalizations of the Theorem of Total Probability. In Sect. 9 we present the founding piece of this work, the Bayes rule for density matrices, its derivation and various properties. We also discuss how the new Bayes rule for density matrices is in some sense the conventional Bayes rule in an optimally chosen eigensystem. Section 10 summarizes all the rules in our calculus and their justifications. In the conclusion section we discuss again how our new calculus relates to quantum physics and possible generalizations of it.

2 Facts on matrices and basic notation

In this paper generalized probability distributions, conditionals and data likelihoods are represented as symmetric positive definite matrices. We will now discuss some relevant matrix algebra facts.

The basic fact that we use a lot is the eigendecomposition of symmetric matrices:

$$\mathbf{S} = \mathbf{S}\boldsymbol{\sigma}\mathbf{S}^\top = \sum_{i=1}^n \sigma_i \mathbf{s}_i \mathbf{s}_i^\top.$$

This says that every such matrix can be written as a product of an orthogonal matrix of eigenvectors \mathbf{S} times a diagonal matrix of eigenvalues $\boldsymbol{\gamma}$ times \mathbf{S}^\top . Alternatively it can be written as mixture of eigendyads formed from the eigenvectors where the eigenvalues act as mixture coefficients.

Any symmetric positive definite³ matrix \mathbf{C} can be seen as a covariance matrix of some random cost vector $\mathbf{c} \in \mathbb{R}^n$, i.e. $\mathbf{C} = \mathbf{E}((\mathbf{c} - \mathbf{E}(\mathbf{c}))(\mathbf{c} - \mathbf{E}(\mathbf{c}))^\top)$. A covariance matrix \mathbf{C} can be depicted as an ellipse $\{\mathbf{C}\mathbf{u} : \|\mathbf{u}\|_2 = 1\}$ centered at the origin, where the eigenvectors form the principal axes and the eigenvalues are the radii of the axes (see Fig. 2).

Note that a covariance matrix \mathbf{C} is diagonal if the components of the cost vector are independent. The variance of the cost vector \mathbf{c} along a vector \mathbf{u} , that is the variance of the dot product $\mathbf{c}^\top \mathbf{u}$, has the form

$$\begin{aligned} \mathbb{V}(\mathbf{c}^\top \mathbf{u}) &= \mathbf{E}((\mathbf{c}^\top \mathbf{u} - \mathbf{E}(\mathbf{c}^\top \mathbf{u}))^2) \\ &= \mathbf{E}(((\mathbf{c}^\top - \mathbf{E}(\mathbf{c}^\top))\mathbf{u})^\top ((\mathbf{c}^\top - \mathbf{E}(\mathbf{c}^\top))\mathbf{u})) \\ &= \mathbf{E}(\mathbf{u}^\top (\mathbf{c} - \mathbf{E}(\mathbf{c}))(\mathbf{c} - \mathbf{E}(\mathbf{c}))^\top \mathbf{u}) \\ &= \mathbf{u}^\top \mathbf{C} \mathbf{u}. \end{aligned}$$

³We use the convention that positive definite matrices have non-negative eigenvalues and *strictly* positive definite matrices have positive eigenvalues.

The variance along an eigenvector of the covariance matrix is the corresponding eigenvalue. Using this interpretation, the matrix \mathbf{C} may be seen as a mapping from the unit ball to $\mathbb{R}_{\geq 0}$, i.e. unit vector \mathbf{u} is mapped to $\mathbf{u}^\top \mathbf{C} \mathbf{u}$. Figure 2 depicts the resulting figure-8-like plots in 2 and 3 dimensions. A second interpretation of the scalar $\mathbf{u}^\top \mathbf{C} \mathbf{u}$ is the square length of \mathbf{u} w.r.t. the basis $\sqrt{\mathbf{C}}$, that is $\mathbf{u}^\top \mathbf{C} \mathbf{u} = \mathbf{u}^\top \sqrt{\mathbf{C}} \sqrt{\mathbf{C}} \mathbf{u} = \|\sqrt{\mathbf{C}} \mathbf{u}\|_2^2$.

The trace $\text{tr}(\mathbf{E})$ of an arbitrary square matrix \mathbf{E} is the sum of its diagonal elements E_{ii} . It is a linear operator. Recall that $\text{tr}(\mathbf{E}\mathbf{F}) = \text{tr}(\mathbf{F}\mathbf{E})$ for any matrices $\mathbf{E} \in \mathbb{R}^{n \times m}$, $\mathbf{F} \in \mathbb{R}^{m \times n}$. Also, for symmetric square matrices, $\text{tr}(\mathbf{S}\mathbf{T}) = \sum_{i,j} S_{ij} T_{ij}$, thus trace can be seen as a dot product between matrices. The trace has a useful cycling property: for arbitrary matrices \mathbf{E} , \mathbf{F} , \mathbf{G} with compatible dimensions $\text{tr}(\mathbf{E}\mathbf{F}\mathbf{G}) = \text{tr}(\mathbf{F}\mathbf{G}\mathbf{E}) = \text{tr}(\mathbf{G}\mathbf{E}\mathbf{F})$. From this follows that trace is *rotation invariant* in the sense that for any orthogonal matrix \mathbf{U} , $\text{tr}(\mathbf{U}\mathbf{E}\mathbf{U}^\top) = \text{tr}(\mathbf{U}^\top \mathbf{U} \mathbf{E}) = \text{tr}(\mathbf{E})$. If \mathbf{S} is symmetric, setting \mathbf{U} to be the eigensystem of \mathbf{S} results in the observation that trace is equal to the sum of eigenvalues of a matrix. Also, for any orthogonal system⁴ \mathbf{u}_i ,

$$\text{tr}(\mathbf{S}) = \text{tr}\left(\underbrace{\sum_{i=1}^n \mathbf{u}_i \mathbf{u}_i^\top}_{\mathbf{I}} \mathbf{S}\right) = \sum_{i=1}^n \mathbf{u}_i^\top \mathbf{S} \mathbf{u}_i.$$

Therefore if \mathbf{S} is symmetric positive definite, then $\text{tr}(\mathbf{S})$ is the total variance along any set of orthogonal directions. Recall that density matrices have trace one and therefore in this case this total variance is always one (see Fig. 3).

The matrix exponential $\exp(\mathbf{S})$ of the symmetric matrix $\mathbf{S} = \sum_i \sigma_i s_i s_i^\top$ is computed by exponentiating the eigenvalues and leaving the eigenvectors unchanged: $\exp(\mathbf{S}) = \sum_i \exp(\sigma_i) s_i s_i^\top$. The matrix logarithm $\log(\mathbf{A})$ is defined similarly but now \mathbf{A} must be strictly positive definite. Clearly, the two functions are inverses of each other. It is important to remember that $\exp(\mathbf{S} + \mathbf{T}) = \exp(\mathbf{S}) \exp(\mathbf{T})$ only holds if \mathbf{S} and \mathbf{T} commute i.e. $\mathbf{S}\mathbf{T} = \mathbf{T}\mathbf{S}$.⁵ However, the following trace inequality, known as the Golden-Thompson inequality⁶ (Bhatia 1997), always holds:

$$\text{tr}(\exp(\mathbf{S}) \exp(\mathbf{T})) \geq \text{tr}(\exp(\mathbf{S} + \mathbf{T})) \quad \text{for symmetric } \mathbf{S} \text{ and } \mathbf{T}, \quad (2.1)$$

where equality holds iff both symmetric matrices commute.

3 Generalized probability distributions and density matrices

In quantum physics a dyad $\mathbf{u}\mathbf{u}^\top$ represents a pure state and density matrices are mixture states. As we shall see density matrices can be interpreted as generalized probability distributions over the set of dyads. Note that in this paper we want to address the statistics community and use linear algebra notation instead of Dirac notation. Any probability vector ($P(M_i)$) can be represented as a diagonal matrix $\text{diag}(P(M_i)) = \sum_i P(M_i) \mathbf{e}_i \mathbf{e}_i^\top$, where \mathbf{e}_i denotes the i th standard basis vector. This means that conventional probability vectors are special density matrices where the eigenvectors are fixed to be the standard basis vectors.

⁴A set of unit vectors \mathbf{u}_i is orthogonal iff $\sum_i \mathbf{u}_i \mathbf{u}_i^\top = \mathbf{I}$.

⁵This occurs iff the two symmetric matrices have the same eigensystem.

⁶Note that the Golden-Thompson inequality does not generalize to three matrices, i.e. there exist symmetric \mathbf{S} , \mathbf{T} , \mathbf{U} , s.t. $\text{tr}(\exp(\mathbf{S}) \exp(\mathbf{T}) \exp(\mathbf{U})) \not\geq \text{tr}(\exp(\mathbf{S} + \mathbf{T} + \mathbf{U}))$.

For the sake of simplicity we assume that our vector space is \mathbb{R}^n . However, everything discussed in this section holds for separable finite or infinite dimensional real and complex Hilbert spaces.

A function $\mu(\mathbf{u})$ from unit vectors \mathbf{u} in \mathbb{R}^n to \mathbb{R} is called a *generalized probability distributions* if the following two conditions hold:

- $\forall \mathbf{u}, 0 \leq \mu(\mathbf{u}) \leq 1$.
- If $\mathbf{u}_1, \dots, \mathbf{u}_n$ form an orthonormal system for \mathbb{R}^n , then $\sum \mu(\mathbf{u}_i) = 1$.

Gleason's Theorem states that there is a one-to-one correspondence between generalized probability distributions and density matrices⁷ in $\mathbb{R}^{n \times n}$:

Theorem 1 (Gleason 1957) *Let $n \geq 3$.⁸ Then any generalized probability distribution μ on \mathbb{R}^n has the form $\mu(\mathbf{u}) = \text{tr}(\mathbf{W}\mathbf{u}\mathbf{u}^\top)$, for a uniquely defined density matrix \mathbf{W} .*

It is easy to see that every density matrix defines a generalized probability distribution. The other direction, is highly non-trivial.⁹ As discussed in the introduction, the dyads $\mathbf{u}\mathbf{u}^\top$ function as elementary events. One may ask what corresponds to arbitrary events and how probabilities are defined for them. In the conventional case, an event is a subset of the domain which can be represented as a vector in $\{0, 1\}^n$. In the generalized setting, an event is a symmetric positive definite matrix \mathbf{P} with eigenvalues in $\{0, 1\}$. Each such matrix \mathbf{P} with eigendecomposition $\sum_{i=1}^k \mathbf{p}_i \mathbf{p}_i^\top$ is a projection matrix for a subspace of \mathbb{R}^n and its probability w.r.t. a distribution \mathbf{W} is defined as the sum of the probabilities of the elementary events $\mathbf{p}_i \mathbf{p}_i^\top$ comprising \mathbf{P} :

$$\text{tr}(\mathbf{W}\mathbf{P}) = \sum_{i=1}^k \text{tr}(\mathbf{W} \mathbf{p}_i \mathbf{p}_i^\top).$$

Interpreting \mathbf{P} as a covariance matrix of some random variable, we can also expand \mathbf{W} and sum the variance along its eigendirections \mathbf{w}_i weighted by the eigenvalues ω_i which are probabilities:

$$\text{tr}(\mathbf{W}\mathbf{P}) = \text{tr}\left(\sum_{i=1}^n \omega_i \mathbf{w}_i \mathbf{w}_i^\top \mathbf{P}\right) = \sum_{i=1}^n \underbrace{\omega_i}_{\text{probability}} \underbrace{\mathbf{w}_i^\top \mathbf{P} \mathbf{w}_i}_{\text{variance}}. \quad (3.1)$$

Random variables are defined in an analogous way. In the conventional case a random variable associates a real value with each point. Now a random variable is an arbitrary symmetric matrix \mathbf{S} . Such matrices have arbitrary real numbers as their eigenvalues and trace $\text{tr}(\mathbf{W}\mathbf{S})$ when \mathbf{S} is expanded becomes the expectation of the random variable w.r.t. den-

⁷The core of the original proof of Gleason's Theorem was for \mathbb{R}^3 (Gleason 1957), and he then generalized the proof to separable real and complex Hilbert spaces of dimension $n \geq 3$.

⁸A slightly different version of this theorem that is based on "effects" instead of dyads holds for dimension 2 as well (Caves et al. 2004).

⁹However, if dyads are replaced by "effects" then the proofs are much simpler (Caves et al. 2004).

sity \mathbf{W} :

$$\text{tr}(\mathbf{W}\mathbf{S}) = \text{tr}\left(\mathbf{W}\sum_i \sigma_i \mathbf{s}_i \mathbf{s}_i^\top\right) = \underbrace{\sum_i \overbrace{\sigma_i}^{\text{outcome}} \overbrace{\mathbf{s}_i^\top \mathbf{W} \mathbf{s}_i}^{\text{probability}}}_{\text{expected outcome}}. \quad (3.2)$$

As discussed before, the conventional case of the expectation calculation is always retained as a special case when all the matrices are diagonal (i.e. fixed eigensystem \mathbf{I}). In quantum physics the expectation calculation $\text{tr}(\mathbf{W}\mathbf{S})$ has the following interpretation: an instrument is represented by a hermitian matrix \mathbf{S} and $\text{tr}(\mathbf{W}\mathbf{S})$ is the expected value of a *quantum measurement* of the mixed state \mathbf{W} with instrument \mathbf{S} . The eigenvalues σ_i of the instrument represent the possible numerical measurement outcomes. Each one of those outcomes is observed with probability $\mathbf{s}_i^\top \mathbf{W} \mathbf{s}_i$, where \mathbf{s}_i is the associated eigenvector of the instrument matrix \mathbf{S} .

In real quantum systems the measurement causes the mixtures state \mathbf{W} to *collapse* into one of the orthogonal states $\{\mathbf{s}_1 \mathbf{s}_1^\top, \dots, \mathbf{s}_n \mathbf{s}_n^\top\}$: the successor state is $\mathbf{s}_i \mathbf{s}_i^\top$ with probability $\mathbf{s}_i^\top \mathbf{W} \mathbf{s}_i$:

$$\mathbf{W} \xrightarrow[\text{collapse}]{\text{measurement}} \underbrace{\sum_i \overbrace{\mathbf{s}_i^\top \mathbf{W} \mathbf{s}_i \mathbf{s}_i \mathbf{s}_i^\top}^{\text{probability}}}_{\text{expected state}}.$$

As we shall see, the expected measurement calculations play an important part in our calculus. However our update rules for density matrices (such as our Bayes rule) do not explicitly include a collapse in the above sense.

Note that some of the equations above hold for arbitrary decompositions into a linear combination of dyads of any size. For example (3.1), holds for any decomposition $\mathbf{W} = \sum_i \omega_i \mathbf{w}_i \mathbf{w}_i^\top$, i.e. the ω_i may be negative, the \mathbf{w}_i may be non-orthogonal, and the size of the decomposition may be larger than n . If the ω_i are non-negative, then they form a probability vector. Similarly, (3.2) also holds for any decomposition $\mathbf{S} = \sum_i \sigma_i \mathbf{s}_i \mathbf{s}_i^\top$. However, quantum measurements are always based on an orthogonal system. Furthermore, orthogonal systems are special in that the orthogonal decomposition of a density matrix $\mathbf{W} = \sum_i \omega_i \mathbf{w}_i \mathbf{w}_i^\top$ attains the minimum of the entropy $\sum_i -\omega_i \ln \omega_i$ over all possible decompositions of \mathbf{W} (inequality (11.86) in Nielsen and Chuang 2000).

A question that naturally arises is whether we can model the generalized probability distributions defined above with a conventional probability space. In other words, is there a conventional probability space and two mappings: one that maps density matrices to conventional probability distributions and the other mapping dyads to events of this probability space. The requirement on these two mappings is that the conventional probability calculations using the images of density matrices and dyads under these mappings satisfy the definition of the generalized probability distributions given above. Essentially, it is known that conventional probability spaces cannot satisfactorily model generalized probabilities, but the details are rather involved. This topic has received considerable attention in the quantum physics community and we refer readers to Holevo (2001) for an extended discussion of impossibility results. Here we only give one simple attempt to model density matrices with a conventional probability space and show that the two natural mappings fail to satisfy the requirements.

A natural interpretation of a density matrix is to view it as a parameterized density over the unit sphere. We claim that if $\mu(\mathbf{u})$ is the uniform density on the sphere, then for any symmetric positive definite matrix $\mathbf{A} \in \mathbb{R}^{n \times n}$ of trace n , $\mathbf{u}^\top \mathbf{A} \mathbf{u} \mu(\mathbf{u})$ is also a conventional probability density on the sphere:

$$\begin{aligned} \int \mathbf{u}^\top \overbrace{\sum_i \alpha_i \mathbf{a}_i \mathbf{a}_i^\top}^{\mathbf{A}} \mathbf{u} \mu(\mathbf{u}) d\mathbf{u} &= \sum_i \alpha_i \int (\mathbf{u}^\top \mathbf{a}_i)^2 \mu(\mathbf{u}) d\mathbf{u} \\ &= \sum_i \alpha_i \int (\mathbf{u}^\top (1, 0, \dots, 0))^2 \mu(\mathbf{u}) d\mathbf{u} \\ &= \text{tr}(\mathbf{A}) \int (u_1)^2 \mu(\mathbf{u}) d\mathbf{u} \\ &= \underbrace{\frac{\text{tr}(\mathbf{A})}{n}}_1 \int \underbrace{\sum_i u_i^2}_1 \mu(\mathbf{u}) d\mathbf{u} = 1. \end{aligned}$$

In the second equality we used the fact that $\mu(\mathbf{u}) d\mathbf{u}$ is uniform and therefore the integral of $(\mathbf{u}^\top \mathbf{a}_i)^2$ is the same as the integral of the squared dot product of \mathbf{u} with the unit vector $(1, 0, \dots, 0)^\top$. In the last equality we again used symmetry: The integral of the square of one component equals the average of the squares of all components.

We modeled density matrices as conventional probability densities over the sphere. Now the natural mapping from dyads to events in the conventional probability space (the sphere) maps $\mathbf{u}\mathbf{u}^\top$ to $\{\mathbf{u}, -\mathbf{u}\}$. However the probability of the latter sets of size 2 is zero with respect to the conventional probabilities densities we defined on the sphere. In particular the probability on any n orthogonal dyads does not sum to one.

4 Commutative matrix product operation

It is well known that the product of two symmetric positive definite matrices might be neither symmetric nor positive definite (see Fig. 4). In this section we define a commutative “product” operation between symmetric positive definite matrices that does result in a symmetric positive definite matrix. Our first definition of this operation requires the two matrices to be strictly positive definite. We then extend the definition to arbitrary symmetric positive definite matrices and prove many properties of this product.

For two symmetric and strictly positive definite matrices \mathbf{A} and \mathbf{B} , we first define the \odot as:

$$\mathbf{A} \odot \mathbf{B} := \exp(\log \overbrace{\mathbf{A} + \mathbf{B}}^{\text{sym.pos.def.}}), \quad (4.1)$$

where here the exponential and logarithm are matrix functions. The matrix log of both matrices produces symmetric matrices which are closed under addition and the matrix exponential of the sum returns a symmetric positive definite matrix. See Fig. 4 for a comparison of matrix product and \odot .

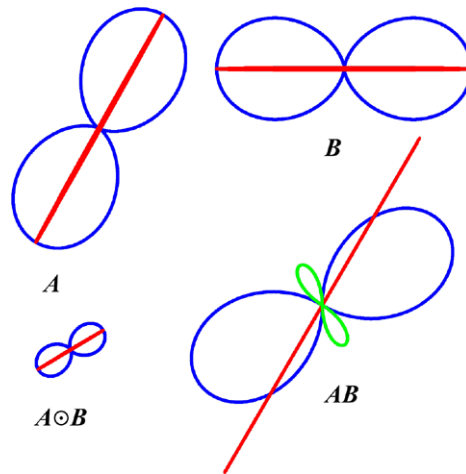


Fig. 4 The matrix product of two positive definite matrices does not preserve positive definiteness. For two matrices A and B we plot their ellipses Au , Bu and figure eights $\text{tr}(Auu^T)u$, $\text{tr}(Buu^T)u$ (for unit u). Both ellipses are very thin, i.e. the ratio between the two eigenvalues of each matrix is 100. We also plot the ellipse ABu and the curve $\text{tr}(ABuu^T)u$. The latter curve consists of two figure eights, the larger one constitutes the part where the trace is positive and the smaller and skinnier one is the part where the trace is negative. This means that AB is not positive definite any more. The product is also not symmetric because the min/max value of $\text{tr}(ABuu^T)$ does not correspond to the axes of the ellipse. Finally, the corresponding plots for $A \odot B$ indicate that this matrix is symmetric and positive definite

Note that we expressed the operation \odot between symmetric strictly positive definite matrices as a $+$ operation between symmetric matrices. Similarly, for any two arbitrary symmetric matrices S and T ,

$$S + T = \log(\exp(S) \odot \exp(T)).$$

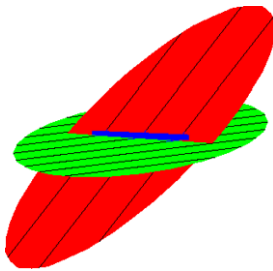
The operation \odot was used in Alexa (2002) to define a “product” between two linear transformations that is commutative. In this paper we use \odot to define conditional density matrices and generalizations of the Bayes rule. A similar path was followed by Cerf and Adami (1999) for defining conditional density matrices of composite systems. We also give a motivation for the operation based on the minimum relative entropy principle (as was done in the conference paper Warmuth 2005) and our probability calculus includes the formula of Cerf and Adami (1999) for composite systems as a special case.

Note that the formula for \odot in (4.1) is not defined if some of the eigenvalues of A or B are zero. We now rewrite the operation using the Lie-Trotter formula and then extend it to arbitrary positive definite matrices. The Lie-Trotter formula (see e.g. Bhatia 1997) is the following equation:

$$\exp(E + F) = \lim_{n \rightarrow \infty} (\exp(E/n) \exp(F/n))^n, \quad \text{any square matrices } E, F.$$

By choosing $E = \log A$ and $F = \log B$, for symmetric and strictly positive definite A and B , we obtain:

$$\exp(\log A + \log B) = \lim_{n \rightarrow \infty} (A^{1/n} B^{1/n})^n.$$



$\text{diag}(\mathbf{A})$	$\text{diag}(\mathbf{B})$	$\text{diag}(\mathbf{AB})$
0	0	0
a	0	0
0	b	0
a	b	ab

Fig. 5 (Color online) When the ellipses \mathbf{A} and \mathbf{B} don't have the same span, then $\mathbf{A} \odot \mathbf{B}$ lies in the intersection of both spans. In the depicted case the intersection is a degenerate ellipse of dimension one (blue line). This generalizes the following intersection property of the matrix product when \mathbf{A} and \mathbf{B} are both diagonal (here of dimension four): $(\mathbf{AB})_{i,i} \neq 0$ iff $\mathbf{A}_{i,i} \neq 0$ and $\mathbf{B}_{i,i} \neq 0$

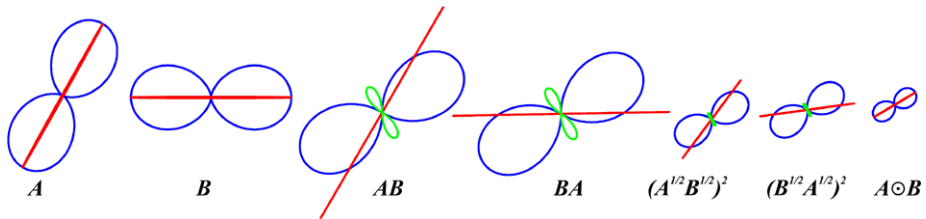


Fig. 6 The behavior of the limit formula for \odot operation. We can see that the additional figure eights indicating negative definiteness are smaller for $(\mathbf{A}^{1/2}\mathbf{B}^{1/2})^2$ than for \mathbf{AB} . As n increases, the additional figure eights shrink further and $\lim_{n \rightarrow \infty} (\mathbf{A}^{1/n}\mathbf{B}^{1/n})^n = \mathbf{A} \odot \mathbf{B}$ becomes positive definite. Also, \mathbf{AB} and \mathbf{BA} are fairly different from one another. The matrices $(\mathbf{A}^{1/2}\mathbf{B}^{1/2})^2$ and $(\mathbf{B}^{1/2}\mathbf{A}^{1/2})^2$ are already more similar and the difference between the two multiplication orders decreases with n until in the limit $\mathbf{A} \odot \mathbf{B} = \mathbf{B} \odot \mathbf{A}$

As n increases, $(\mathbf{A}^{1/n}\mathbf{B}^{1/n})^n$ gets closer and closer to being positive definite and symmetric. The first couple iterations of the limit formula are plotted in Fig. 6. See Alexa (2002) for additional plots. Notice that the limit is defined even when \mathbf{A} and \mathbf{B} have zero eigenvalues. We therefore extend the definition of \odot to arbitrary symmetric positive definite matrices \mathbf{A} and \mathbf{B} :

$$\mathbf{A} \odot \mathbf{B} := \lim_{n \rightarrow \infty} (\mathbf{A}^{1/n}\mathbf{B}^{1/n})^n. \quad (4.2)$$

From now on we use the above extended definition of \odot . Numerous properties of this operation are given below.

Theorem 2 For any symmetric positive definite matrices $\mathbf{A}, \mathbf{B}, \mathbf{C}$ the following holds:

OP1 Intersection property:

$$\text{range}(\mathbf{A} \odot \mathbf{B}) = \text{range}(\mathbf{A}) \cap \text{range}(\mathbf{B}).$$

where the range of a matrix is the linear subspace spanned by the columns of the matrix. This property generalizes the intersection properties for products of diagonal matrices (which model conventional probability distributions): the product of two diagonal matrices with the characteristic vectors of two subsets as diagonals gives a diagonal matrix formed from the characteristic vector of the intersection (see Fig. 5).

OP2 Let \mathbf{R}_A be a matrix whose columns form an orthonormal basis for the range of \mathbf{A} , i.e. $\mathbf{R}_A \in \mathbb{R}^{n \times k}$ and $\mathbf{R}_A^\top \mathbf{R}_A = \mathbf{I}_k$, where k is the dimensionality of the range of \mathbf{A} . Define \mathbf{R}_B analogously. In a similar fashion $\mathbf{R}_{A \cap B}$ will contain the basis for the intersection of ranges. Let \log^+ denote the modified matrix logarithm that takes the log of non-zero eigenvalues but leaves the zero eigenvalues unchanged. This operation can be also defined by the following formula:¹⁰

$$\log^+ \mathbf{A} = \mathbf{R}_A \log(\mathbf{R}_A^\top \mathbf{A} \mathbf{R}_A) \mathbf{R}_A^\top. \quad (4.3)$$

With this notation, \odot can be written as

$$\mathbf{A} \odot \mathbf{B} = \mathbf{R}_{A \cap B} \exp(\mathbf{R}_{A \cap B}^\top (\log^+ \mathbf{A} + \log^+ \mathbf{B}) \mathbf{R}_{A \cap B}) \mathbf{R}_{A \cap B}^\top. \quad (4.4)$$

OP3 $\mathbf{A} \odot \mathbf{B} = \mathbf{AB}$ if \mathbf{A} and \mathbf{B} commute.

OP4 \odot is commutative, i.e. $\mathbf{A} \odot \mathbf{B} = \mathbf{B} \odot \mathbf{A}$.

OP5 The identity matrix is the neutral element, i.e. $\mathbf{A} \odot \mathbf{I} = \mathbf{A}$.

OP6 $(c\mathbf{A}) \odot \mathbf{B} = c(\mathbf{A} \odot \mathbf{B})$, for any scalar $c > 0$.

OP7 $\mathbf{A} \odot \mathbf{A}^{-1} = \mathbf{I}$ for invertible \mathbf{A} . Also, $\mathbf{A} \odot \mathbf{A}^+ = \mathbf{P}_A$, where \mathbf{A}^+ denotes the pseudoinverse and \mathbf{P}_A is the projection matrix¹¹ for range(\mathbf{A}).

OP8 \odot is associative, i.e. $(\mathbf{A} \odot \mathbf{B}) \odot \mathbf{C} = \mathbf{A} \odot (\mathbf{B} \odot \mathbf{C})$.

OP9 Monotonic convergence of the limit defining \odot :

$$\forall n \geq 1 : \text{tr}(\mathbf{A}^{1/(n+1)} \mathbf{B}^{1/(n+1)})^{n+1} \leq \text{tr}(\mathbf{A}^{1/n} \mathbf{B}^{1/n})^n.$$

OP10 $\text{tr}(\mathbf{A} \odot \mathbf{B}) \leq \text{tr}(\mathbf{AB})$, where equality holds iff \mathbf{A} and \mathbf{B} commute. In particular, for any unit \mathbf{u} $\text{tr}(\mathbf{A} \odot \mathbf{uu}^\top) = \text{tr}(\mathbf{Auu}^\top)$ iff \mathbf{u} is an eigenvector of \mathbf{A} .

OP11 For any unit direction $\mathbf{u} \in \text{range}(\mathbf{A})$, $\mathbf{A} \odot \mathbf{uu}^\top = e^{\mathbf{u}^\top (\log^+ \mathbf{A}) \mathbf{u}} \mathbf{uu}^\top$.

OP12 For any unit direction \mathbf{u} and eigendecomposition $\sum_i \alpha_i \mathbf{a}_i \mathbf{a}_i^\top$ of a strictly positive definite matrix \mathbf{A} ,

$$\text{tr}(\mathbf{Auu}^\top) = \sum_i (\mathbf{u}^\top \mathbf{a}_i)^2 \alpha_i, \quad \text{and} \quad \text{tr}(\mathbf{A} \odot \mathbf{uu}^\top) = \prod_i \alpha_i^{(\mathbf{u}^\top \mathbf{a}_i)^2},$$

i.e. the matrix product corresponds to an arithmetic average and the \odot product to a geometric average of the eigenvalues of \mathbf{A} .

OP13 $\det(\mathbf{A} \odot \mathbf{B}) = \det(\mathbf{A}) \det(\mathbf{B})$, which is the same as for the normal matrix product.

OP14 For any orthogonal system \mathbf{u}_i , we have $\prod_i \text{tr}(\mathbf{A} \odot \mathbf{u}_i \mathbf{u}_i^\top) = \det(\mathbf{A})$.

OP15 For any unit direction \mathbf{u} , $\text{tr}((\mathbf{A} \odot \mathbf{B}) \odot \mathbf{uu}^\top) = \text{tr}(\mathbf{A} \odot \mathbf{uu}^\top) \text{tr}(\mathbf{B} \odot \mathbf{uu}^\top)$.

OP16 For any unit direction $\mathbf{u} \in \text{range} \mathbf{A}$, $\text{tr}(\mathbf{A}^+ \odot \mathbf{uu}^\top) = \frac{1}{\text{tr}(\mathbf{A} \odot \mathbf{uu}^\top)}$, where \mathbf{A}^+ denotes the pseudoinverse.

Proof Properties OP1 and OP2 follow from results in Kato (1978) or Theorem 1.2 of Simon (1979). Here we only prove that $\text{range}(\mathbf{A} \odot \mathbf{B}) \subseteq \text{range}(\mathbf{A}) \cap \text{range}(\mathbf{B})$. We can split the limit

¹⁰Note that when the rank k of \mathbf{A} is zero, then one still can define the projections in a consistent manner. In this case \mathbf{R}_A is of dimension $n \times 0$, and the matrices $\mathbf{R}_A^\top \mathbf{R}_A$ and $\log(\mathbf{R}_A^\top \mathbf{E} \mathbf{R}_A)$ are of dimension 0×0 for any $\mathbf{E} \in \mathbb{R}^{n \times n}$. Also it is natural to define $\mathbf{R}_A \mathbf{E} \mathbf{R}_A^\top$ as the $n \times n$ zero matrix $\mathbf{0}$. With this definition, the r.h.s. of (4.3) is $\mathbf{0}$ when \mathbf{A} is $\mathbf{0}$.

¹¹Note that $\mathbf{P}_A = \mathbf{R}_A \mathbf{R}_A^\top$.

defining \odot as follows:

$$A \odot B = \lim_{n \rightarrow \infty} (A^{1/n} B^{1/n})^n = \lim_{n \rightarrow \infty} A^{1/n} \lim_{n \rightarrow \infty} B^{1/n} (A^{1/n} B^{1/n})^{n-1}. \quad (4.5)$$

Here we used the property that $\lim E_n F_n = \lim E_n \lim F_n$ if all the limits exist. This follows from the corresponding sum and product properties of scalar limits and the fact that entries of a product matrix are finite sums of products.

It is easy to see that $\lim_{n \rightarrow \infty} A^{1/n} = P_A$ because the matrix power for a symmetric matrix corresponds to taking powers of the eigenvalues and n -th roots converge to either zero or one. Thus the limit is a matrix whose eigenvalues are 0 or 1, which is a projection matrix. By plugging

$$\lim_{n \rightarrow \infty} A^{1/n} = P_A = P_A P_A = P_A \lim_{n \rightarrow \infty} A^{1/n}$$

into (4.5) we get

$$A \odot B = P_A \lim_{n \rightarrow \infty} A^{1/n} \lim_{n \rightarrow \infty} B^{1/n} (A^{1/n} B^{1/n})^{n-1} = P_A (A \odot B).$$

This implies that $\text{range}(A \odot B) \subseteq \text{range}(A)$. Similarly we can prove $A \odot B = (A \odot B) P_B$, which implies that $\text{range}(A \odot B) \subseteq \text{range}(B)$, and therefore $\text{range}(A \odot B) \subseteq \text{range}(A) \cap \text{range}(B)$.

Property OP3 can be seen from the definition of \odot via the limit formula (4.2): when A and B commute, then the n copies of $A^{1/n}$ in $(A^{1/n} B^{1/n})^n$ can be gathered into A and similarly for $B^{1/n}$.

Properties 4–7 easily follow from the formula (4.4) for \odot .

Property OP8. For strictly positive definite A, B, C associativity reduces to the associativity of addition in the log domain. To show it in general we use the representation (4.4) of \odot via the \log^+ operation. Let $R = R_{(A \cap B) \cap C} = R_{A \cap (B \cap C)}$. Then:

$$(A \odot B) \odot C = R \exp(R^\top (\log^+(A \odot B) + \log^+ C) R) R^\top.$$

Now we rewrite $\log^+(A \odot B)$ using (4.3):

$$\log^+(A \odot B) = R_{A \cap B} \log(R_{A \cap B}^\top (A \odot B) R_{A \cap B}) R_{A \cap B}^\top.$$

Substituting expression (4.4) for $A \odot B$ into the above and using $R_{A \cap B}^\top R_{A \cap B} = I_k$ we get:

$$\log^+(A \odot B) = \underbrace{R_{A \cap B} R_{A \cap B}^\top}_{P_{A \cap B}} (\log^+ A + \log^+ B) \underbrace{R_{A \cap B} R_{A \cap B}^\top}_{P_{A \cap B}}. \quad (4.6)$$

Here $P_{A \cap B} = R_{A \cap B} R_{A \cap B}^\top$ is the projection matrix onto the subspace $\text{range}(A) \cap \text{range}(B)$. All the basis vectors of $A \cap B \cap C$ obviously lie in the larger subspace as well, thus the projection leaves them unchanged and we get $P_{A \cap B} R = R$, $R^\top P_{A \cap B} = R^\top$. Thus:

$$(A \odot B) \odot C = R \exp(R^\top (\log^+ A + \log^+ B + \log^+ C) R) R^\top.$$

The same expression can be obtained for $A \odot (B \odot C)$, thus establishing associativity.

Property OP9. By Fact 8.10.9 of Bernstein (2005), we have that for any positive definite matrices A and B , and $n \geq 1$:

$$\text{tr}(A^n B^n)^{n+1} \leq \text{tr}(A^{n+1} B^{n+1})^n.$$

Now by substituting $A = A^{1/(n(n+1))}$, $B = B^{1/(n(n+1))}$ the monotonicity property OP9 immediately follows.

Property OP10. When A and B are strictly positive definite, this inequality is an instantiation of the Golden-Thompson inequality (2.1). For arbitrary positive definite matrices, the property follows from the previous monotonicity property OP9. Note that there are symmetric positive definite matrices A , B and C s.t. $\text{tr}(A \odot B \odot C) \not\leq \text{tr}(ABC)$.

Property OP11. We use the expression for \odot operation given in (4.4). Since $u \in \text{range}(A)$, the basis of the intersection space is u itself:

$$uu^\top \odot A = u \exp(u^\top (\log^+ uu^\top + \log^+ A) u) u^\top.$$

Note that $\log^+ uu^\top = 0$ and that the expression inside the exponential is a scalar. The desired property immediately follows by moving this scalar to the front.

Property OP12. The expression for the trace of the matrix product is the expected measurement interpretation (3.2) discussed in Sect. 3. Note that $(u^\top a_i)^2$ is a probability vector and in this expression uu^\top can be replaced by any density matrix.

For the second trace $\text{tr}(A \odot uu^\top)$, we can rewrite it using OP11 and eigendecomposition of A as follows:

$$\text{tr}(A \odot uu^\top) = e^{u^\top \log A u} = e^{\sum_i (a_i \cdot u)^2 \log \alpha_i} = \prod_i \alpha_i^{(u^\top a_i)^2},$$

which is a weighted geometric average of α_i with weights $(u^\top a_i)^2$.

Property OP13. Since $\det(EF) = \det(E)\det(F)$, and for symmetric matrices S and $r \in \mathbb{R}$, $\det(S^r) = \det(S)^r$, we have $\det((A^{1/n} B^{1/n})^n) = \det(A)\det(B)$ for all $n \in \mathbb{N}$. By property (4.2), the limit of the l.h.s. of the last equality becomes $A \odot B$ and this proves the property.

Property OP14. If A is not full rank, then $\det(A)$ is zero. In that case, there will be some u_i that is not in the range of A . For that u_i , $\text{tr}(A \odot u_i u_i^\top) = 0$, making the whole product zero. When A has full rank, we rewrite the product as follows:

$$\begin{aligned} \prod_i \text{tr}(A \odot u_i u_i^\top) &\stackrel{\text{OP11}}{=} \prod_i e^{\text{tr}(\log A u_i u_i^\top)} \\ &= e^{\text{tr}(\log A \sum u_i u_i^\top)} = e^{\text{tr}(\log A)} = \prod_i \alpha_i = \det(A). \end{aligned}$$

Property OP15. If $u \notin \text{range}(A) \cap \text{range}(B) \stackrel{\text{OP1}}{=} \text{range}(A \odot B)$, then the property trivially holds because $\text{tr}((A \odot B) \odot uu^\top)$ and either $\text{tr}(A \odot uu^\top)$ or $\text{tr}(B \odot uu^\top)$ are zero. When $u \in \text{range}(A) \cap \text{range}(B)$, then the property essentially follows from $e^{a+b} = e^a e^b$:

$$\begin{aligned} \text{tr}((A \odot B) \odot uu^\top) &\stackrel{\text{OP11}}{=} e^{u^\top \log^+(A \odot B) u} \stackrel{(4.6)}{=} e^{u^\top P_{A \cap B} (\log^+ A + \log^+ B) P_{A \cap B} u} = e^{u^\top (\log^+ A + \log^+ B) u} \\ &= e^{u^\top \log^+ A u} e^{u^\top \log^+ B u} = \text{tr}(A \odot uu^\top) \text{tr}(B \odot uu^\top). \end{aligned}$$

Needless to say Property OP15 does not hold if uu^\top is replaced by a mixture of dyads.

Property OP16. Trivially follows from OP11. \square

We will now discuss some of the properties further. In particular, we will show a simple example that demonstrates that the upper bound OP10 can be quite loose when both matrices

are dyads. In this case the inequality becomes:

$$\text{tr}(\mathbf{u}\mathbf{u}^\top \odot \mathbf{v}\mathbf{v}^\top) \leq \text{tr}(\mathbf{u}\mathbf{u}^\top \mathbf{v}\mathbf{v}^\top) = (\mathbf{u} \cdot \mathbf{v})^2.$$

The right hand side can be made arbitrarily close to one by choosing almost parallel \mathbf{u} and \mathbf{v} . The left side is zero in this case, which can be seen by analyzing the intersection of the ranges. Dyads are rank one matrices and their ranges are lines through the origin. The intersection of two such lines is either only the origin or the line itself. Thus, by Property OP1 it follows that $\mathbf{u}\mathbf{u}^\top \odot \mathbf{v}\mathbf{v}^\top = \mathbf{0}$, unless $\mathbf{u} = \pm\mathbf{v}$. This can also be seen from the limit expression in (4.2):

$$\begin{aligned} \mathbf{u}\mathbf{u}^\top \odot \mathbf{v}\mathbf{v}^\top &= \lim_{n \rightarrow \infty} ((\mathbf{u}\mathbf{u}^\top)^{\frac{1}{n}} (\mathbf{v}\mathbf{v}^\top)^{\frac{1}{n}})^n = \lim_{n \rightarrow \infty} (\mathbf{u}\mathbf{u}^\top \mathbf{v}\mathbf{v}^\top)^n \\ &= \left(\lim_{n \rightarrow \infty} (\mathbf{u} \cdot \mathbf{v})^{2n-1} \right) \mathbf{u}\mathbf{v}^\top = \mathbf{0}, \quad \text{unless } \mathbf{u} = \pm\mathbf{v}. \end{aligned}$$

Where the last equality holds because $|\mathbf{u} \cdot \mathbf{v}| < 1$, when $\mathbf{u} \neq \pm\mathbf{v}$.

Note that the expression (4.4) for \odot based on \log^+ gives us a convenient method for computing the operation even when the matrices have some zero eigenvalues. The modified matrix logarithm \log^+ is easily computed via (4.3). The matrix \mathbf{R}_A containing the orthonormal basis for range of \mathbf{A} can be computed using Gram-Schmidt orthogonalization procedure or the QR-decomposition. To compute the basis for the intersection of $\text{range}(\mathbf{A})$ and $\text{range}(\mathbf{B})$, we express the intersection i.t.o. the union and the orthogonal complement $^\perp$ of a space:

$$\text{range}(\mathbf{A}) \cap \text{range}(\mathbf{B}) = (\text{range}(\mathbf{A})^\perp \cup \text{range}(\mathbf{B})^\perp)^\perp.$$

For any matrix \mathbf{E} , an orthonormal basis for $\text{range}(\mathbf{E})^\perp$ can be obtained by completing an orthonormal basis for $\text{range}(\mathbf{E})$ to an orthonormal basis for the whole space. The additional basis vectors needed are the basis for $\text{range}(\mathbf{E})^\perp$. Also, if we have two matrices \mathbf{E} and \mathbf{F} , we can get the range for the union of their ranges just by putting all columns of \mathbf{E} and \mathbf{F} together into a bigger matrix $\mathbf{G} = (\mathbf{E}, \mathbf{F})$. Clearly, $\text{range}(\mathbf{G}) = \text{range}(\mathbf{E}) \cup \text{range}(\mathbf{F})$. Piecing all of this together gives an implementation of the \odot operation.

5 Joint distributions

A density matrix defines a generalized probability distribution over the dyads from one space. However we need to consider several spaces and joint distributions over them. In the conventional case A, B denote finite sets $\{a_1, \dots, a_{n_A}\}, \{b_1, \dots, b_{n_B}\}$, $(P(a_i)), (P(b_j))$ probability vectors over these sets and $(P(a_i, b_j))$ is an $n_A \times n_B$ dimensional matrix of probabilities for the tuple set $A \times B$. In the generalized case, \mathbb{A}, \mathbb{B} denote real finite dimensional vector spaces of dimension $n_{\mathbb{A}}, n_{\mathbb{B}}$ and $\mathbf{D}(\mathbb{A}), \mathbf{D}(\mathbb{B})$ are the density matrices defining the generalized probability distributions over these spaces. The *joint space* (\mathbb{A}, \mathbb{B}) is the tensor product¹² between the spaces \mathbb{A} and \mathbb{B} , which is of dimension $n_{\mathbb{A}}n_{\mathbb{B}}$. The joint distribution is specified by a density matrix over this joint space, denoted by $\mathbf{D}(\mathbb{A}, \mathbb{B})$.

¹²See Bhatia (1997) for a formal definition of tensor product between vector spaces. For us, the tensor product of $\mathbb{R}^{n_{\mathbb{A}}}$ and $\mathbb{R}^{n_{\mathbb{B}}}$ is $\mathbb{R}^{n_{\mathbb{A}}n_{\mathbb{B}}}$.

We let $\mathbf{D}(\mathbf{a})$, $\mathbf{D}(\mathbf{b})$ denote the probabilities assigned to dyads $\mathbf{a}\mathbf{a}^\top$, $\mathbf{b}\mathbf{b}^\top$ from the spaces \mathbb{A} , \mathbb{B} by the density matrices $\mathbf{D}(\mathbb{A})$, $\mathbf{D}(\mathbb{B})$, respectively:

$$\mathbf{D}(\mathbf{a}) := \text{tr}(\mathbf{D}(\mathbb{A})\mathbf{a}\mathbf{a}^\top), \quad \mathbf{D}(\mathbf{b}) := \text{tr}(\mathbf{D}(\mathbb{B})\mathbf{b}\mathbf{b}^\top). \quad (\text{MJ1})$$

The conventional probability distributions can be seen as diagonal density matrices. A probability distribution $(P(a_i))$ on the set A is the density matrix $\text{diag}((P(a_i)))$. Also $P(a_j) = \mathbf{e}_j^\top \text{diag}((P(a_i))) \mathbf{e}_j$.

To introduce the joint probability $\mathbf{D}(\mathbf{a}, \mathbf{b})$ we need the Kronecker matrix product. Given two matrices \mathbf{E} and \mathbf{F} with dimensions $n \times m$ and $p \times q$, their Kronecker product (also known as direct product or tensor product) $\mathbf{E} \otimes \mathbf{F}$ is a matrix with dimensions $np \times mq$ which in block form is given as:

$$\begin{pmatrix} e_{11}\mathbf{F} & e_{12}\mathbf{F} & \dots & e_{1m}\mathbf{F} \\ e_{21}\mathbf{F} & e_{22}\mathbf{F} & \dots & e_{2m}\mathbf{F} \\ \dots & \dots & \dots & \dots \\ e_{n1}\mathbf{F} & e_{n2}\mathbf{F} & \dots & e_{nm}\mathbf{F} \end{pmatrix}.$$

The Kronecker product has the following useful properties:

KP1 $(\mathbf{E} \otimes \mathbf{F})^\top = \mathbf{E}^\top \otimes \mathbf{F}^\top$.

KP2 $(\mathbf{E} \otimes \mathbf{F})(\mathbf{G} \otimes \mathbf{H}) = \mathbf{E}\mathbf{G} \otimes \mathbf{F}\mathbf{H}$ if the dimensions are appropriate.

KP3 $\text{tr}(\mathbf{E} \otimes \mathbf{F}) = \text{tr}(\mathbf{E})\text{tr}(\mathbf{F})$.

KP4 If symmetric matrix \mathbf{S} has eigenvalues σ_i and eigenvectors \mathbf{s}_i and symmetric matrix \mathbf{T} has eigenvalues τ_j and eigenvectors \mathbf{t}_j , then $\mathbf{S} \otimes \mathbf{T}$ has eigenvalues $\sigma_i \tau_j$ and eigenvectors $\mathbf{s}_i \otimes \mathbf{t}_j$.

KP5 For symmetric positive definite matrices $\mathbf{A}, \mathbf{B}, \mathbf{C}, \mathbf{D}$, $(\mathbf{A} \otimes \mathbf{B}) \odot (\mathbf{C} \otimes \mathbf{D}) = (\mathbf{A} \odot \mathbf{C}) \otimes (\mathbf{B} \odot \mathbf{D})$.

The first four properties are standard. The last property follows from the limit definition (4.2) of the \odot operation.

$$\begin{aligned} (\mathbf{A} \otimes \mathbf{B}) \odot (\mathbf{C} \otimes \mathbf{D}) &= \lim_{n \rightarrow \infty} ((\mathbf{A} \otimes \mathbf{B})^{\frac{1}{n}} (\mathbf{C} \otimes \mathbf{D})^{\frac{1}{n}})^n \\ &= \lim_{n \rightarrow \infty} ((\mathbf{A}^{\frac{1}{n}} \mathbf{C}^{\frac{1}{n}}) \otimes (\mathbf{B}^{\frac{1}{n}} \mathbf{D}^{\frac{1}{n}}))^n \\ &= \left(\lim_{n \rightarrow \infty} (\mathbf{A}^{\frac{1}{n}} \mathbf{B}^{\frac{1}{n}})^n \right) \otimes \left(\lim_{n \rightarrow \infty} (\mathbf{C}^{\frac{1}{n}} \mathbf{D}^{\frac{1}{n}})^n \right). \end{aligned}$$

The last transition which moved the limit inside the Kronecker product, follows from the fact that the elements of the Kronecker product matrix are just pairwise products of elements from the two matrices. And when all limits exist, a limit of a product of two number sequences is a product of limits.

Now the joint probability $\mathbf{D}(\mathbf{a}, \mathbf{b})$ becomes the probability assigned by density matrix $\mathbf{D}(\mathbb{A}, \mathbb{B})$ to the jointly specified dyad $(\mathbf{a} \otimes \mathbf{b})(\mathbf{a} \otimes \mathbf{b})^\top$:

$$\mathbf{D}(\mathbf{a}, \mathbf{b}) := \text{tr}(\mathbf{D}(\mathbb{A}, \mathbb{B})(\mathbf{a} \otimes \mathbf{b})(\mathbf{a} \otimes \mathbf{b})^\top) = \text{tr}(\mathbf{D}(\mathbb{A}, \mathbb{B})(\mathbf{a}\mathbf{a}^\top \otimes \mathbf{b}\mathbf{b}^\top)). \quad (\text{MJ3})$$

Note that in the conventional case a joint probability between two sets A and B is defined over all pairs of points from A and B . This corresponds to the so-called “separable case”

in the generalized setting when the joint density $D(\mathbb{A}, \mathbb{B})$ can be expressed as $\sum_i \sigma_i \mathbf{a}_i \mathbf{a}_i^\top \otimes \mathbf{b}_i \mathbf{b}_i^\top$, where the \mathbf{a}_i and \mathbf{b}_i are states of spaces \mathbb{A} and \mathbb{B} , respectively. However, in the more general case, there are elementary events in the joint space (\mathbb{A}, \mathbb{B}) that don't decompose into elementary events of the spaces \mathbb{A} and \mathbb{B} , i.e. there are dyads in the joint space that are not of the form $\mathbf{a}\mathbf{a}^\top \otimes \mathbf{b}\mathbf{b}^\top$, where $\mathbf{a}\mathbf{a}^\top$ and $\mathbf{b}\mathbf{b}^\top$ are dyads of \mathbb{A} and \mathbb{B} , respectively. In quantum physics the inseparability of a density matrix is called “entanglement”.

6 Marginalization of the joint via partial traces

We would like to be able to perform marginalization operations on our joint density matrix $D(\mathbb{A}, \mathbb{B})$, i.e. obtain the density matrix $D(\mathbb{A})$ from the joint matrix. In the conventional case the marginalization was performed by summing out one of the variables by summing the rows or the columns of the matrix specifying the joint probability distribution. For density matrices, the analogous operation is the *partial trace* (see e.g. Nielsen and Chuang 2000).

The partial trace is a generalization of normal matrix trace. It typically produces a matrix instead of a number and can be used to retrieve the (scaled) factor matrices from a Kronecker product. We denote the partial trace with $\text{tr}_{\mathbb{A}}$, where \mathbb{A} specifies the space to be “summed out”. Suppose G is a matrix over the space $\mathbb{A} \otimes \mathbb{B}$ and \mathbb{A} has dimension n and \mathbb{B} dimension m . Thus G has dimension $nm \times nm$ and can be written in block form as a $n \times n$ matrix of $m \times m$ matrices G_{ij} :

$$G = \begin{pmatrix} G_{11} & G_{12} & \dots & G_{1n} \\ G_{21} & G_{22} & \dots & G_{2n} \\ \dots & \dots & \dots & \dots \\ G_{n1} & G_{n2} & \dots & G_{nn} \end{pmatrix}.$$

Here we suppose that space \mathbb{A} is \mathbb{R}^n and space \mathbb{B} is \mathbb{R}^m . Then the two partial traces of this matrix are given by:

$$\underbrace{\text{tr}_{\mathbb{A}}(G)}_{m \times m} = G_{11} + G_{22} + \dots + G_{nn},$$

$$\underbrace{\text{tr}_{\mathbb{B}}(G)}_{n \times n} = \begin{pmatrix} \text{tr}(G_{11}) & \text{tr}(G_{12}) & \dots & \text{tr}(G_{1n}) \\ \text{tr}(G_{21}) & \text{tr}(G_{22}) & \dots & \text{tr}(G_{2n}) \\ \dots & \dots & \dots & \dots \\ \text{tr}(G_{n1}) & \text{tr}(G_{n2}) & \dots & \text{tr}(G_{nn}) \end{pmatrix}.$$

In multilinear algebra partial traces are known as tensor contractions and can of course be generalized to the tensor product of more than two spaces. The partial trace is a linear operator and we now give some other useful properties:

$$\text{PT1 } \text{tr}_{\mathbb{A}}(E \otimes F) = \text{tr}(E)F, \quad \text{tr}_{\mathbb{B}}(E \otimes F) = \text{tr}(F)E.$$

$$\text{PT2 } \text{tr}(G) = \text{tr}(\text{tr}_{\mathbb{A}}(G)) = \text{tr}(\text{tr}_{\mathbb{B}}(G)).$$

$$\text{PT3 } \text{tr}_{\mathbb{A}}(G(I_{\mathbb{A}} \otimes F)) = \text{tr}_{\mathbb{A}}(G)F, \quad \text{tr}_{\mathbb{A}}((I_{\mathbb{A}} \otimes F)G) = F\text{tr}_{\mathbb{A}}(G).$$

$$\text{PT4 } \text{tr}(G(E \otimes F)) = \text{tr}(\text{tr}_{\mathbb{B}}(G(I_{\mathbb{A}} \otimes F))E).$$

The first three properties are straightforward and the last one follows from the others as follows:

$$\begin{aligned}\mathrm{tr}(\mathbf{G}(\mathbf{E} \otimes \mathbf{F})) &\stackrel{\text{KP2}}{=} \mathrm{tr}(\mathbf{G}(\mathbf{I}_{\mathbb{A}} \otimes \mathbf{F})(\mathbf{E} \otimes \mathbf{I}_{\mathbb{B}})) = \mathrm{tr}((\mathbf{E} \otimes \mathbf{I}_{\mathbb{B}})\mathbf{G}(\mathbf{I}_{\mathbb{A}} \otimes \mathbf{F})) \\ &\stackrel{\text{PT2}}{=} \mathrm{tr}(\mathrm{tr}_{\mathbb{B}}((\mathbf{E} \otimes \mathbf{I}_{\mathbb{B}})\mathbf{G}(\mathbf{I}_{\mathbb{A}} \otimes \mathbf{F}))) \stackrel{\text{PT3}}{=} \mathrm{tr}(\mathbf{E}\mathrm{tr}_{\mathbb{B}}(\mathbf{G}(\mathbf{I}_{\mathbb{A}} \otimes \mathbf{F}))).\end{aligned}$$

We use the partial trace to define marginals as follows:

$$\mathbf{D}(\mathbb{A}) := \mathrm{tr}_{\mathbb{B}}(\mathbf{D}(\mathbb{A}, \mathbb{B})), \quad \mathbf{D}(\mathbb{B}) := \mathrm{tr}_{\mathbb{A}}(\mathbf{D}(\mathbb{A}, \mathbb{B})). \quad (\text{MJ2})$$

The following lemma shows that $\mathbf{D}(\mathbb{A})$ and $\mathbf{D}(\mathbb{B})$ defined this way are again density matrices.

Lemma 1 *Partial trace of a density matrix is also a density matrix.*

Proof Symmetry is obvious. Trace one follows from Property PT2 of the partial trace:

$$\mathrm{tr}(\mathbf{D}(\mathbb{A})) = \mathrm{tr}(\mathrm{tr}_{\mathbb{B}}(\mathbf{D}(\mathbb{A}, \mathbb{B}))) = \mathrm{tr}(\mathbf{D}(\mathbb{A}, \mathbb{B})) = 1.$$

Positive definiteness follows by a similar argument:

$$\begin{aligned}\mathbf{a}^{\top} \mathbf{D}(\mathbb{A}) \mathbf{a} &= \mathrm{tr}(\mathbf{D}(\mathbb{A}) \mathbf{a} \mathbf{a}^{\top}) \stackrel{\text{PT3}}{=} \mathrm{tr}(\mathrm{tr}_{\mathbb{B}}(\mathbf{D}(\mathbb{A}, \mathbb{B})(\mathbf{a} \mathbf{a}^{\top} \otimes \mathbf{I}_{\mathbb{B}}))) \\ &\stackrel{\text{PT2}}{=} \mathrm{tr}\left(\mathbf{D}(\mathbb{A}, \mathbb{B}) \underbrace{\left(\mathbf{a} \mathbf{a}^{\top} \otimes \sum_i \mathbf{b}_i \mathbf{b}_i^{\top}\right)}_{\mathbf{I}_{\mathbb{B}}}\right) \\ &= \sum_i \mathrm{tr}(\mathbf{D}(\mathbb{A}, \mathbb{B})(\mathbf{a} \mathbf{a}^{\top} \otimes \mathbf{b}_i \mathbf{b}_i^{\top})) \\ &= \sum_i (\mathbf{a} \otimes \mathbf{b}_i)^{\top} \mathbf{D}(\mathbb{A}, \mathbb{B})(\mathbf{a} \otimes \mathbf{b}_i) \geq 0.\end{aligned}$$

□

Partial traces also allow us to define objects of the type $\mathbf{D}(\mathbb{A}, \mathbf{b})$. In the conventional case this corresponds to taking one row or column out of the joint probability table. In the generalized case we want the following property to be satisfied:

$$\mathrm{tr}(\mathbf{D}(\mathbb{A}, \mathbf{b}) \mathbf{a} \mathbf{a}^{\top}) = \mathbf{D}(\mathbf{a}, \mathbf{b}). \quad (\text{MJ5})$$

This is accomplished by defining $\mathbf{D}(\mathbb{A}, \mathbf{b})$ via the following formula:

$$\mathbf{D}(\mathbb{A}, \mathbf{b}) := \mathrm{tr}_{\mathbb{B}}(\mathbf{D}(\mathbb{A}, \mathbb{B})(\mathbf{I}_{\mathbb{A}} \otimes \mathbf{b} \mathbf{b}^{\top})). \quad (\text{MJ4})$$

Property MJ5 now follows from partial trace Property PT4. We can also see that trace of $\mathbf{D}(\mathbf{A}, \mathbf{b})$ gives us the probability $\mathbf{D}(\mathbf{b})$:

$$\begin{aligned} \text{tr}(\mathbf{D}(\mathbb{A}, \mathbf{b})) &\stackrel{\text{MJ4}}{=} \text{tr}(\text{tr}_{\mathbb{B}}(\mathbf{D}(\mathbb{A}, \mathbb{B})(\mathbf{I}_{\mathbb{A}} \otimes \mathbf{b}\mathbf{b}^{\top}))) \stackrel{\text{PT2}}{=} \text{tr}(\text{tr}_{\mathbb{A}}(\mathbf{D}(\mathbb{A}, \mathbb{B})(\mathbf{I}_{\mathbb{A}} \otimes \mathbf{b}\mathbf{b}^{\top}))) \\ &\stackrel{\text{PT3}}{=} \text{tr}(\text{tr}_{\mathbb{A}}(\mathbf{D}(\mathbb{A}, \mathbb{B}))\mathbf{b}\mathbf{b}^{\top}) \stackrel{\text{MJ2}}{=} \text{tr}(\mathbf{D}(\mathbb{B})\mathbf{b}\mathbf{b}^{\top}) = \mathbf{D}(\mathbf{b}). \end{aligned} \quad (6.1)$$

A brief note on matrix properties of $\mathbf{D}(\mathbb{A}, \mathbf{b})$. We just saw that its trace is $\mathbf{D}(\mathbf{b})$ which is between zero and one. Since it satisfies Property (MJ5), it is positive definite as well. Symmetry is also easily verified.

Note that for any orthogonal system \mathbf{b}_i of \mathbb{B} ,

$$\mathbf{D}(\mathbb{A}) = \sum_i \mathbf{D}(\mathbb{A}, \mathbf{b}_i).$$

This can be seen as follows.

$$\begin{aligned} \mathbf{D}(\mathbb{A}) &\stackrel{\text{MJ2}}{=} \text{tr}_{\mathbb{B}}(\mathbf{D}(\mathbb{A}, \mathbb{B})) = \text{tr}_{\mathbb{B}}(\mathbf{D}(\mathbb{A}, \mathbb{B})(\mathbf{I}_{\mathbb{A}} \otimes \mathbf{I}_{\mathbb{B}})) \\ &= \text{tr}_{\mathbb{B}}\left(\mathbf{D}(\mathbb{A}, \mathbb{B})\left(\mathbf{I}_{\mathbb{A}} \otimes \underbrace{\sum_i \mathbf{b}_i \mathbf{b}_i^{\top}}_{\mathbf{I}_{\mathbb{B}}}\right)\right) = \sum_i \text{tr}_{\mathbb{B}}(\mathbf{D}(\mathbb{A}, \mathbb{B})(\mathbf{I}_{\mathbb{A}} \otimes \mathbf{b}_i \mathbf{b}_i^{\top})) \\ &\stackrel{\text{MJ4}}{=} \sum_i \mathbf{D}(\mathbb{A}, \mathbf{b}_i). \end{aligned}$$

The conventional definition of independence also naturally generalizes: $\mathbf{D}(\mathbb{A})$ is independent of $\mathbf{D}(\mathbb{B})$ if the joint density matrix decomposes: $\mathbf{D}(\mathbb{A}, \mathbb{B}) = \mathbf{D}(\mathbb{A}) \otimes \mathbf{D}(\mathbb{B})$. It is easy to see that in this case we have $\mathbf{D}(\mathbf{a}, \mathbf{b}) = \mathbf{D}(\mathbf{a})\mathbf{D}(\mathbf{b})$ for all \mathbf{a}, \mathbf{b} :

$$\begin{aligned} \mathbf{D}(\mathbf{a}, \mathbf{b}) &= \text{tr}((\mathbf{D}(\mathbb{A}) \otimes \mathbf{D}(\mathbb{B}))(\mathbf{a}\mathbf{a}^{\top} \otimes \mathbf{b}\mathbf{b}^{\top})) \stackrel{\text{KP2}}{=} \text{tr}((\mathbf{D}(\mathbb{A})\mathbf{a}\mathbf{a}^{\top}) \otimes (\mathbf{D}(\mathbb{B})\mathbf{b}\mathbf{b}^{\top})) \\ &\stackrel{\text{KP3}}{=} \text{tr}(\mathbf{D}(\mathbb{A})\mathbf{a}\mathbf{a}^{\top})\text{tr}(\mathbf{D}(\mathbb{B})\mathbf{b}\mathbf{b}^{\top}) = \mathbf{D}(\mathbf{a})\mathbf{D}(\mathbf{b}). \end{aligned}$$

7 Conditional probabilities

The topic of conditional probabilities in this generalized setting contains many subtleties. First we will give the defining formulas for conditional density matrices and then discuss some of the issues.

CP1 $\mathbf{D}(\mathbb{A}|\mathbb{B}) := \mathbf{D}(\mathbb{A}, \mathbb{B}) \odot (\mathbf{I}_{\mathbb{A}} \otimes \mathbf{D}(\mathbb{B}))^{-1}$ (formula (4) of Cerf and Adami 1999 expressed with the \odot operation). This formula requires $\mathbf{D}(\mathbb{B})$ to be invertible. In the conventional case, this corresponds to the conditional probabilities being undefined if the event conditioned on has probability zero.

CP2 $\mathbf{D}(\mathbb{A}|\mathbf{b}) := \frac{\mathbf{D}(\mathbb{A}, \mathbf{b})}{\mathbf{D}(\mathbf{b})}$.

CP3 $\mathbf{D}(\mathbf{a}|\mathbb{B}) := \mathbf{D}(\mathbf{a}, \mathbb{B}) \odot \mathbf{D}(\mathbb{B})^{-1}$.

CP4 $\mathbf{D}(\mathbf{a}|\mathbf{b}) := \frac{\mathbf{D}(\mathbf{a}, \mathbf{b})}{\mathbf{D}(\mathbf{b})}$. This basic conditional probability is a straightforward generalization of the conventional case. It also has a quantum-mechanical interpretation. See Appendix for details.

Note that CP1 has the form: density matrix \odot inverse of a normalization. We can also reexpress the other definitions in this unified form:

$$\begin{aligned}
\text{CP'2 } D(\mathbb{A}|\mathbb{B}) &= \underbrace{\text{tr}_{\mathbb{B}}(D(\mathbb{A}, \mathbb{B})(I_{\mathbb{A}} \otimes \mathbf{b}\mathbf{b}^{\top}))}_{D(\mathbb{A}, \mathbf{b})} \odot \underbrace{\text{tr}_{\mathbb{B}}((I_{\mathbb{A}} \otimes D(\mathbb{B}))(I_{\mathbb{A}} \otimes \mathbf{b}\mathbf{b}^{\top}))^{-1}}_{\frac{1}{D(\mathbf{b})} I_{\mathbb{A}}} \\
\text{CP'3 } D(\mathbf{a}|\mathbb{B}) &= \underbrace{\text{tr}_{\mathbb{A}}(D(\mathbb{A}, \mathbb{B})(\mathbf{a}\mathbf{a}^{\top} \otimes I_{\mathbb{B}}))}_{D(\mathbf{a}, \mathbb{B})} \odot \underbrace{\text{tr}_{\mathbb{A}}((I_{\mathbb{A}} \otimes D(\mathbb{B}))(\mathbf{a}\mathbf{a}^{\top} \otimes I_{\mathbb{B}}))^{-1}}_{D(\mathbb{B})^{-1}} \\
\text{CP'4 } D(\mathbf{a}|\mathbf{b}) &= \underbrace{\text{tr}(D(\mathbb{A}, \mathbb{B})(\mathbf{a}\mathbf{a}^{\top} \otimes \mathbf{b}\mathbf{b}^{\top}))}_{D(\mathbf{a}, \mathbf{b})} \odot \underbrace{\text{tr}((I_{\mathbb{A}} \otimes D(\mathbb{B}))(\mathbf{a}\mathbf{a}^{\top} \otimes \mathbf{b}\mathbf{b}^{\top}))^{-1}}_{\frac{1}{D(\mathbf{b})}}
\end{aligned}$$

We say that the joint density $D(\mathbb{A}, \mathbb{B})$ is *decoupled* if its eigendecomposition has the form: $D(\mathbb{A}, \mathbb{B}) = (\mathcal{W}_{\mathbb{A}} \otimes \mathcal{W}_{\mathbb{B}}) \text{diag}(\omega)(\mathcal{W}_{\mathbb{A}} \otimes \mathcal{W}_{\mathbb{B}})^{\top}$. Note that $\mathcal{W}_{\mathbb{A}} \otimes \mathcal{W}_{\mathbb{B}}$ is orthogonal iff both $\mathcal{W}_{\mathbb{A}}$ and $\mathcal{W}_{\mathbb{B}}$ are orthogonal. As we shall see later, dealing with conditionals is often simpler in the decoupled case. We first prove an upper bound for $\text{tr}(D(\mathbb{A}|\mathbb{B}))$ that is tight iff the joint is decoupled.

Lemma 2 *The following inequality holds:*

$$\text{tr}(D(\mathbb{A}|\mathbb{B})) \leq n_{\mathbb{B}},$$

where $n_{\mathbb{B}}$ is the dimensionality of space \mathbb{B} . Furthermore, $\text{tr}(D(\mathbb{A}|\mathbb{B})) = n_{\mathbb{B}}$ if and only if the joint $D(\mathbb{A}, \mathbb{B})$ is decoupled.

Proof The inequality is shown using properties of \odot and partial traces:

$$\begin{aligned}
\text{tr}(D(\mathbb{A}|\mathbb{B})) &\stackrel{\text{CP1}}{=} \text{tr}(D(\mathbb{A}, \mathbb{B}) \odot (I_{\mathbb{A}} \otimes D(\mathbb{B})^{-1})) \stackrel{\text{OP10}}{\leq} \text{tr}(D(\mathbb{A}, \mathbb{B})(I_{\mathbb{A}} \otimes D(\mathbb{B})^{-1})) \\
&\stackrel{\text{PT2}}{=} \text{tr}(\text{tr}_{\mathbb{A}}(D(\mathbb{A}, \mathbb{B})(I_{\mathbb{A}} \otimes D(\mathbb{B})^{-1}))) \stackrel{\text{PT3}}{=} \text{tr}(\underbrace{\text{tr}_{\mathbb{A}}(D(\mathbb{A}, \mathbb{B}))}_{D(\mathbb{B})} D(\mathbb{B})^{-1}) \\
&= \text{tr}(I_{\mathbb{B}}) = n_{\mathbb{B}}.
\end{aligned}$$

Remember that equality in Property OP10 of \odot only occurs when the two matrices commute. Two matrices commute iff their eigensystems are the same. This gives us the condition that the eigensystem of $D(\mathbb{A}, \mathbb{B})$ must be the same as the eigensystem of $I_{\mathbb{A}} \otimes D(\mathbb{B})^{-1}$. The latter eigensystem is clearly decoupled. Thus for equality to hold it is *necessary* that the eigensystem of $D(\mathbb{A}, \mathbb{B})$ be decoupled.

Now we will argue that it is also *sufficient*. Let the joint density matrix have eigensystem $D(\mathbb{A}, \mathbb{B}) = (\mathcal{W}_{\mathbb{A}} \otimes \mathcal{W}_{\mathbb{B}}) \text{diag}(\omega)(\mathcal{W}_{\mathbb{A}}^{\top} \otimes \mathcal{W}_{\mathbb{B}}^{\top})$. $I_{\mathbb{A}}$ commutes with any matrix on space \mathbb{A} . Therefore it suffices to show that the marginal $D(\mathbb{B})$ in this case has eigensystem $\mathcal{W}_{\mathbb{B}}$. The decoupled eigensystem matrix $\mathcal{W}_{\mathbb{A}} \otimes \mathcal{W}_{\mathbb{B}}$ has the following list of $n_{\mathbb{A}}n_{\mathbb{B}}$ columns:

$$\begin{aligned}
\mathcal{W}_{\mathbb{A}} \otimes \mathcal{W}_{\mathbb{B}} &= \left(\mathbf{w}_{\mathbb{A}}^1 \otimes \mathbf{w}_{\mathbb{B}}^1, \mathbf{w}_{\mathbb{A}}^1 \otimes \mathbf{w}_{\mathbb{B}}^2, \dots, \mathbf{w}_{\mathbb{A}}^1 \otimes \mathbf{w}_{\mathbb{B}}^{n_{\mathbb{B}}}, \right. \\
&\quad \mathbf{w}_{\mathbb{A}}^2 \otimes \mathbf{w}_{\mathbb{B}}^1, \mathbf{w}_{\mathbb{A}}^2 \otimes \mathbf{w}_{\mathbb{B}}^2, \dots, \mathbf{w}_{\mathbb{A}}^2 \otimes \mathbf{w}_{\mathbb{B}}^{n_{\mathbb{B}}}, \\
&\quad \dots, \dots, \dots, \dots, \dots, \dots, \\
&\quad \left. \mathbf{w}_{\mathbb{A}}^{n_{\mathbb{A}}} \otimes \mathbf{w}_{\mathbb{B}}^1, \mathbf{w}_{\mathbb{A}}^{n_{\mathbb{A}}} \otimes \mathbf{w}_{\mathbb{B}}^2, \dots, \mathbf{w}_{\mathbb{A}}^{n_{\mathbb{A}}} \otimes \mathbf{w}_{\mathbb{B}}^{n_{\mathbb{B}}} \right).
\end{aligned}$$

In correspondence with this structure we adopt a double indexing scheme for the eigenvalues $\omega_{i,j}$ of the joint matrix $D(\mathbb{A}, \mathbb{B})$, where $\omega_{i,j}$ is the eigenvalue associated with eigenvector

$\mathbf{w}_{\mathbb{A}}^i \otimes \mathbf{w}_{\mathbb{B}}^j$. The index i runs from 1 to $n_{\mathbb{A}}$, and j runs to $n_{\mathbb{B}}$. Now the eigendecomposition can be written as:

$$\mathbf{D}(\mathbb{A}, \mathbb{B}) = \sum_{i,j} \omega_{i,j} (\mathbf{w}_{\mathbb{A}}^i (\mathbf{w}_{\mathbb{A}}^i)^{\top} \otimes \mathbf{w}_{\mathbb{B}}^j (\mathbf{w}_{\mathbb{B}}^j)^{\top}).$$

Partial trace is a linear operator and $\text{tr}_{\mathbb{A}}(\mathbf{w}_{\mathbb{A}}^i (\mathbf{w}_{\mathbb{A}}^i)^{\top} \otimes \mathbf{w}_{\mathbb{B}}^j (\mathbf{w}_{\mathbb{B}}^j)^{\top}) \stackrel{\text{PT1}}{=} \mathbf{w}_{\mathbb{B}}^j (\mathbf{w}_{\mathbb{B}}^j)^{\top}$. Therefore:

$$\mathbf{D}(\mathbb{B}) = \text{tr}_{\mathbb{A}}(\mathbf{D}(\mathbb{A}, \mathbb{B})) = \sum_j \omega_j \mathbf{w}_{\mathbb{B}}^j (\mathbf{w}_{\mathbb{B}}^j)^{\top},$$

where $\omega_j = \sum_i \omega_{i,j}$. Thus we produced the eigendecomposition of the marginal $\mathbf{D}(\mathbb{B})$ and it indeed has eigensystem $\mathcal{W}_{\mathbb{B}}$. \square

Let us briefly discuss the connection and difference between our notion of decoupled joints and the notion of entanglement that appears in quantum physics. Recall that entanglement, as we mentioned at the end of Sect. 5 corresponds to the fact that there are dyads $\mathbf{c}\mathbf{c}^{\top}$ in the joint space (\mathbb{A}, \mathbb{B}) that can't be written as $\mathbf{a}\mathbf{a}^{\top} \otimes \mathbf{b}\mathbf{b}^{\top}$ for any two dyads $\mathbf{a}\mathbf{a}^{\top}$ and $\mathbf{b}\mathbf{b}^{\top}$ in \mathbb{A} and \mathbb{B} . This notion carries over to mixed states or density matrices. In quantum physics, a joint density matrix $\mathbf{D}(\mathbb{A}, \mathbb{B})$ is called *separable* (or non-entangled) if it can be expressed as $(\mathcal{W}_{\mathbb{A}} \otimes \mathcal{W}_{\mathbb{B}}) \text{diag}(\omega) (\mathcal{W}_{\mathbb{A}} \otimes \mathcal{W}_{\mathbb{B}})^{\top}$. The crucial difference between the definitions of separable and decoupled matrices is that in the separable case, $\mathcal{W}_{\mathbb{A}}$ and $\mathcal{W}_{\mathbb{B}}$ don't have to be orthogonal. Every decoupled matrix is separable, but there are separable density matrices that are not decoupled. The question of deciding whether a given matrix is separable is known to be very difficult, whereas the question of being decoupled is easily decided by e.g. the condition of the above lemma. One of the reasons for which Cerf and Adami (1999) introduced a conditional density matrix via Rule CP1 was to give a necessary condition for the separability of a joint density matrix.

To complete the rules for conditional density matrices, we would need rules that allow us to marginalize the conditionals, e.g. for going from $\mathbf{D}(\mathbb{A}|\mathbb{B})$ to $\mathbf{D}(\mathbf{a}|\mathbf{b})$. One obvious consequence of our definitions is the marginalization rule for $\mathbf{D}(\mathbb{A}|\mathbf{b})$:

$$\mathbf{D}(\mathbf{a}|\mathbf{b}) \stackrel{\text{CP4}}{=} \frac{\mathbf{D}(\mathbf{a}, \mathbf{b})}{\mathbf{D}(\mathbf{b})} \stackrel{\text{MJ5}}{=} \frac{\text{tr}(\mathbf{D}(\mathbb{A}, \mathbf{b}) \mathbf{a}\mathbf{a}^{\top})}{\mathbf{D}(\mathbf{b})} \stackrel{\text{CP2}}{=} \text{tr}(\mathbf{D}(\mathbb{A}|\mathbf{b}) \mathbf{a}\mathbf{a}^{\top}). \quad (\text{MC4})$$

There don't seem to be any other simple marginalization rules for $\mathbf{D}(\mathbb{A}|\mathbb{B})$ and $\mathbf{D}(\mathbf{a}|\mathbb{B})$ that hold for arbitrary joints. However, when the joint is decoupled, then the following additional marginalization rule for $\mathbf{D}(\mathbb{A}|\mathbb{B})$ is valid:

Lemma 3 *For all decoupled joints $\mathbf{D}(\mathbb{A}, \mathbb{B})$,*

$$\mathbf{D}(\mathbf{a}|\mathbb{B}) = \text{tr}_{\mathbb{A}}(\mathbf{D}(\mathbb{A}|\mathbb{B})(\mathbf{a}\mathbf{a}^{\top} \otimes \mathbf{I}_{\mathbb{B}})).$$

Proof We will compute both sides of the equation and show them to be identical. We begin by writing down the decomposition of the decoupled joint from Lemma 2:

$$\mathbf{D}(\mathbb{A}, \mathbb{B}) = \sum_{i,j} \omega_{i,j} (\mathbf{w}_{\mathbb{A}}^i (\mathbf{w}_{\mathbb{A}}^i)^{\top} \otimes \mathbf{w}_{\mathbb{B}}^j (\mathbf{w}_{\mathbb{B}}^j)^{\top}). \quad (7.1)$$

Additionally, in the same lemma, the following form for $\mathbf{D}(\mathbb{B})$ was established in this case:

$$\mathbf{D}(\mathbb{B}) = \text{tr}_{\mathbb{A}}(\mathbf{D}(\mathbb{A}, \mathbb{B})) = \sum_j \omega_j \mathbf{w}_{\mathbb{B}}^j (\mathbf{w}_{\mathbb{B}}^j)^{\top}, \quad (7.2)$$

where $\omega_j = \sum_i \omega_{i,j}$. According to CP3, $D(a|\mathbb{B}) = D(a, \mathbb{B}) \odot D(\mathbb{B})^{-1}$, therefore we will need to compute $D(a, \mathbb{B})$:

$$\begin{aligned} D(a, \mathbb{B}) &\stackrel{\text{MJ4}}{=} \text{tr}_{\mathbb{A}}(D(\mathbb{A}, \mathbb{B})(aa^{\top} \otimes I_{\mathbb{B}})) \stackrel{(7.1)}{=} \sum_{i,j} \omega_{i,j} \text{tr}_{\mathbb{A}}((w_{\mathbb{A}}^i (w_{\mathbb{A}}^i)^{\top} \otimes w_{\mathbb{B}}^j (w_{\mathbb{B}}^j)^{\top})(aa^{\top} \otimes I_{\mathbb{B}})) \\ &\stackrel{\text{KP2}}{=} \sum_{i,j} \omega_{i,j} \text{tr}_{\mathbb{A}}((w_{\mathbb{A}}^i (w_{\mathbb{A}}^i)^{\top}) aa^{\top} \otimes w_{\mathbb{B}}^j (w_{\mathbb{B}}^j)^{\top}) \stackrel{\text{PT1}}{=} \sum_{i,j} \omega_{i,j} (w_{\mathbb{A}}^i \cdot a)^2 w_{\mathbb{B}}^j (w_{\mathbb{B}}^j)^{\top}. \end{aligned}$$

Together with (7.2), this gives:

$$D(a|\mathbb{B}) = \sum_{i,j} \frac{\omega_{i,j} (w_{\mathbb{A}}^i \cdot a)^2}{\omega_j} w_{\mathbb{B}}^j (w_{\mathbb{B}}^j)^{\top}.$$

Now, we proceed to the right side of the equation in the lemma. Substituting (7.1) and (7.2) into the formula for $D(\mathbb{A}|\mathbb{B})$ we obtain:

$$D(\mathbb{A}|\mathbb{B}) \stackrel{\text{CP1}}{=} D(\mathbb{A}, \mathbb{B}) \odot (I_{\mathbb{A}} \otimes D(\mathbb{B})^{-1}) = \sum_{i,j} \frac{\omega_{i,j}}{\omega_j} (w_{\mathbb{A}}^i (w_{\mathbb{A}}^i)^{\top} \otimes w_{\mathbb{B}}^j (w_{\mathbb{B}}^j)^{\top}). \quad (7.3)$$

Using linearity of the partial trace we compute the right side as follows:

$$\begin{aligned} \text{tr}_{\mathbb{A}}(D(\mathbb{A}|\mathbb{B})(aa^{\top} \otimes I_{\mathbb{B}})) &\stackrel{(7.3)}{=} \sum_{i,j} \frac{\omega_{i,j}}{\omega_j} \text{tr}_{\mathbb{A}}((w_{\mathbb{A}}^i (w_{\mathbb{A}}^i)^{\top} \otimes w_{\mathbb{B}}^j (w_{\mathbb{B}}^j)^{\top})(aa^{\top} \otimes I_{\mathbb{B}})) \\ &\stackrel{\text{KP2}}{=} \sum_{i,j} \frac{\omega_{i,j}}{\omega_j} \text{tr}_{\mathbb{A}}((w_{\mathbb{A}}^i (w_{\mathbb{A}}^i)^{\top}) aa^{\top} \otimes w_{\mathbb{B}}^j (w_{\mathbb{B}}^j)^{\top}) \\ &\stackrel{\text{PT1}}{=} \sum_{i,j} \frac{\omega_{i,j} (w_{\mathbb{A}}^i \cdot a)^2}{\omega_j} w_{\mathbb{B}}^j (w_{\mathbb{B}}^j)^{\top}. \end{aligned} \quad \square$$

As discussed, obtaining $D(a|b) = \frac{\text{tr}(D(\mathbb{A}, \mathbb{B})(aa^{\top} \otimes bb^{\top}))}{\text{tr}(D(\mathbb{B})bb^{\top})}$ from $D(\mathbb{A}|\mathbb{B})$ is non-trivial. In particular, there are cases where

$$\text{tr}(D(\mathbb{A}|\mathbb{B})(aa^{\top} \otimes bb^{\top})) \neq D(a|b),$$

even when $D(\mathbb{A}, \mathbb{B})$ is decoupled and a and b are not eigenvectors of $D(\mathbb{A})$ and $D(\mathbb{B})$, respectively. Curiously enough, if we replace the matrix product with \odot , then we always have

$$\begin{aligned} &\text{tr}(D(\mathbb{A}|\mathbb{B}) \odot (aa^{\top} \otimes bb^{\top})) \\ &\stackrel{\text{CP1}}{=} \text{tr}((D(\mathbb{A}, \mathbb{B}) \odot (I_{\mathbb{A}} \otimes D(\mathbb{B})^{-1})) \odot (aa^{\top} \otimes bb^{\top})) \\ &\stackrel{\text{OP15}}{=} \text{tr}(D(\mathbb{A}, \mathbb{B}) \odot (aa^{\top} \otimes bb^{\top})) \text{tr}((I_{\mathbb{A}} \otimes D(\mathbb{B})^{-1}) \odot (aa^{\top} \otimes bb^{\top})) \\ &\stackrel{\text{OP16}}{=} \frac{\text{tr}(D(\mathbb{A}, \mathbb{B}) \odot (aa^{\top} \otimes bb^{\top}))}{\text{tr}(D(\mathbb{B}) \odot bb^{\top})}. \end{aligned}$$

Let us now recall the conditionals in the conventional probability theory. The full conditional table $P(A|B)$ lists conditional probabilities of all pairs of elementary events $P(a_i|b_j)$.

This table has the obvious properties: The sum of all entries is n_B and the sum of any column is 1, i.e. $\sum_i P(a_i|b_j) = \sum_i \frac{P(a_i, b_j)}{P(b_j)} = \frac{P(b_j)}{P(b_j)} = 1$. Thus a conditional table is a column-stochastic matrix and for any such matrix we can construct a joint that has that matrix as its conditional table. For example we can take arbitrary probability vector p and multiply the i -th column of $P(A|B)$ by p_i , now the sum of each column is p_i and thus the sum of all entries is 1 and we have a valid joint. Note that this implies that many different joints have the same conditional table.

The decoupled case behaves as the conventional case, i.e. many joints correspond to the same conditional. A decoupled joint and conditional always have the same eigensystem and going from the joint to the conditional is similar to the conventional case (see (7.1–7.3) for details).

However, for non-decoupled joint density matrices, i.e. when $\text{tr}(D(A|B)) < n_B$ (Lemma 2), the situation is quite different. For example, the eigenvalues of $D(A|B)$ can now be bigger than 1 (Cerf and Adami 1999). Also based on numerical experiments, we conjecture that in the non-decoupled case, the mapping between $D(A, B)$ and $D(A|B)$ is invertible, i.e. unlike the conventional case, there is only one joint that gives rise to a given conditional matrix. In other words we conjecture that in the non-decoupled case it suffices to specify the conditional $D(A|B)$.

More specifically, we claim that the following EM-like algorithm converges to $D(B)$ and then $D(A, B) \stackrel{\text{CP1}}{=} D(A|B) \odot D(B)$: W_0 is initialized to I_B/n_B and the estimate W_{t+1} for $D(B)$ is computed from $D(A|B)$ and the previous estimate W_t as

$$W_{t+1} = \frac{\text{tr}_A(D(A|B) \odot (I_A \otimes W_t))}{\text{tr}(D(A|B) \odot (I_A \otimes W_t))}.$$

8 Theorems of total probability

The Theorem of Total Probability is an important calculation in conventional probability theory. It expresses probability of some event a as an expected conditional probability of the elementary events b_i that form a partition of the probability space B :

$$P(a) = \sum_i P(a|b_i)P(b_i).$$

TP1 For any orthogonal system b_i of B , $D(a) = \sum_i D(a|b_i)D(b_i)$.

TP2 $D(a) = \text{tr}(D(a|B) \odot D(B))$.

TP3 $D(A) = \text{tr}_B(D(A|B) \odot (I_A \otimes D(B)))$.

The first formula can be shown as follows:

$$D(a) = \text{tr}(D(a, B)) = \sum_i b_i^\top D(a, B) b_i = \sum_i D(a, b_i) = \sum_i D(a|b_i)D(b_i).$$

To derive the second apply $\odot D(B)$ to both sides of CP3, take trace of both sides and use (6.1). The proof of the third property follows the same outline but uses CP1 and MJ2.

Conventional versions of the last two properties are obtained when the density and conditional matrices are diagonal. Note that in general these generalizations of the Theorem of Total Probability do not “decouple”, i.e. you cannot write them as a sum of products of conditional and marginal probabilities. However, using the Property OP10 of \odot operation we

can establish upper bounds on probability of $\mathbf{D}(\mathbf{a})$ in terms of “decoupled” sums that look like the conventional versions of the Theorem of Total Probability. If $\mathbf{D}(\mathbb{B}) = \sum_i \omega_i \mathbf{w}_i \mathbf{w}_i^\top$ and $\mathbf{D}(\mathbf{a}|\mathbb{B}) = \sum_i \lambda_i \mathbf{u}_i \mathbf{u}_i^\top$ are eigendecompositions of the corresponding matrices, then

$$\begin{aligned}
 \mathbf{D}(\mathbf{a}) &= \text{tr}(\mathbf{D}(\mathbf{a}|\mathbb{B}) \odot \mathbf{D}(\mathbb{B})) \leq \text{tr}(\mathbf{D}(\mathbf{a}|\mathbb{B}) \mathbf{D}(\mathbb{B})) \\
 &= \underbrace{\sum_i \underbrace{\omega_i}_{\text{probability}} \underbrace{\mathbf{w}_i^\top \mathbf{D}(\mathbf{a}|\mathbb{B}) \mathbf{w}_i}_{\text{variance}}}_{\text{expected variance}} \\
 &= \underbrace{\sum_i \underbrace{\mathbf{u}_i^\top \mathbf{D}(\mathbb{B}) \mathbf{u}_i}_{\text{probability } \mathbf{D}(\mathbf{u}_i)} \underbrace{\lambda_i}_{\text{outcome}}}_{\text{expected measurement}}. \tag{8.1}
 \end{aligned}$$

The first version of the upper bound corresponds to using the eigendecomposition of $\mathbf{D}(\mathbb{B})$ and can be interpreted as an expected variance calculation with $\mathbf{D}(\mathbf{a}|\mathbb{B})$ as the covariance matrix. The second version expands $\mathbf{D}(\mathbf{a}|\mathbb{B})$ and corresponds to a quantum measurement of system in state $\mathbf{D}(\mathbb{B})$ with instrument specified by $\mathbf{D}(\mathbf{a}|\mathbb{B})$. Letting $p(b_i)$ equal ω_i or $\mathbf{u}_i^\top \mathbf{D}(\mathbb{B}) \mathbf{u}_i$ and letting $p(a_i|b_i)$ equal $\mathbf{w}_i^\top \mathbf{D}(\mathbf{a}|\mathbb{B}) \mathbf{w}_i$ or λ_i , we see the correspondence of these upper bounds to the conventional Theorem of Total Probability. The equality only occurs when $\mathbf{D}(\mathbf{a}|\mathbb{B})$ and $\mathbf{D}(\mathbb{B})$ commute.

9 Bayes rules

In the conventional setup we assume that a model M_i is chosen with prior probability $P(M_i)$. The model then generates the data y with probability $P(y|M_i)$, i.e.

$$\begin{aligned}
 P(y) &= \sum_i P(M_i) P(y|M_i) \\
 &= \text{tr}(\text{diag}((P(M_i))) \text{diag}((P(y|M_i))).
 \end{aligned}$$

The reason why we expressed $P(y)$ as a trace of two diagonal matrices will become apparent in a moment.

The generalized setup is completely analogous. There is an underlying joint space (\mathbb{M}, \mathbb{Y}) between the model space \mathbb{M} and the data space \mathbb{Y} . The prior is specified by a density matrix $\mathbf{D}(\mathbb{M})$. The data is a unit direction \mathbf{y} in \mathbb{Y} space that is generated by the density $\mathbf{D}(\mathbb{Y})$. The probability $\mathbf{D}(\mathbf{y})$ can be expressed i.t.o. the prior $\mathbf{D}(\mathbb{M})$ and data likelihood $\mathbf{D}(\mathbf{y}|\mathbb{M})$ using TP2:

$$\mathbf{D}(\mathbf{y}) = \text{tr}(\mathbf{D}(\mathbb{M}) \odot \mathbf{D}(\mathbf{y}|\mathbb{M})).$$

Note that in the conventional case we first chose a model based on the prior and then generated data based on the chosen model. In the generalized case we do not know how to decouple the action on the prior from the choice of the data when conditioned on the prior.

Let us first recall the conventional Bayes rule and rewrite it in matrix notation:

$$P(M_i|y) = \frac{P(M_i)P(y|M_i)}{P(y)}, \quad \text{where } P(y) = \sum_j P(M_j)P(y|M_j), \quad (9.1)$$

$$\text{diag}(P(M_i|y)) = \frac{\text{diag}(P(M_i)) \text{diag}(P(y|M_i))}{\text{tr}(\text{diag}(P(M_i)) \text{diag}(P(y|M_i)))}.$$

We now present and discuss the analogous Bayes rule for the generalized setting. At the end of this section we present a list of all Bayes rules.

In the generalized Bayes rule we cannot simply multiply the prior density matrix with the data likelihood matrix. This is because a product of two symmetric positive definite matrices can be neither symmetric nor positive definite (see Fig. 4). Instead, we replace the matrix multiplication with \odot operation:

$$D(\mathbb{M}|y) = \frac{D(\mathbb{M}) \odot D(y|\mathbb{M})}{D(y)}, \quad \text{where } D(y) = \text{tr}(D(\mathbb{M}) \odot D(y|\mathbb{M})). \quad (9.2)$$

Normalizing by the trace ensures that the trace of the posterior density matrix is one. In both the conventional as well as the new Bayes rule above, the normalization constant is the likelihood of the data. When the matrices $D(\mathbb{M})$ and $D(y|\mathbb{M})$ have the same eigensystem, then \odot becomes the matrix multiplication. In the following subsections we derive the above Bayes rules from the minimum relative entropy principle. For the conventional Bayes rule the standard relative entropy between probability vectors is used, whereas the generalized Bayes rule and the crucial \odot operation is motivated by the quantum relative entropy between density matrices due to Umegaki (see e.g. Nielsen and Chuang 2000).

We visualize the conventional Bayes rule in Fig. 7. Repeated application of the rule with the same likelihood makes the posteriors increasingly concentrated on the point with maximum data likelihood $P(y|M_i)$. Therefore this rule can be interpreted as a soft maximum-likelihood calculation. Figure 8 demonstrates the generalized Bayes rule. There the posterior gradually moves towards the eigenvector belonging to the largest eigenvalue of the data likelihood matrix $D(y|\mathbb{M})$. Thus the new rule can be interpreted as a soft calculation of the eigenvector with maximum eigenvalue.

In Fig. 11 we depict a sequence of updates with the new Bayes rule when the data likelihood matrix is different in each iteration. Observe that based on the relative lengths of the axes (eigenvalues) and the directions of the axes (eigenvectors) in the ellipse describing the current data likelihood matrix, the posterior adjusts its axis lengths and directions.

Other Bayes rules for our calculus are listed below. They all express one conditional in terms of the corresponding reverse conditional.

$$\text{BR1 } D(\mathbb{B}|\mathbb{A}) = (I_{\mathbb{A}} \otimes D(\mathbb{B})) \odot D(\mathbb{A}|\mathbb{B}) \odot (D(\mathbb{A}) \otimes I_{\mathbb{B}})^{-1}, \quad \text{where } D(\mathbb{A}) = \text{tr}_{\mathbb{B}}((I_{\mathbb{A}} \otimes D(\mathbb{B})) \odot D(\mathbb{A}|\mathbb{B})).$$

$$\text{BR2 } D(\mathbb{b}|\mathbb{A}) = D(\mathbb{b})D(\mathbb{A}|\mathbb{b}) \odot D(\mathbb{A})^{-1}, \quad \text{where } D(\mathbb{A}) \stackrel{\text{TP3}}{=} \text{tr}_{\mathbb{B}}((I_{\mathbb{A}} \otimes D(\mathbb{B})) \odot D(\mathbb{A}|\mathbb{B})).$$

$$\text{BR3 } D(\mathbb{B}|\mathbf{a}) = \frac{D(\mathbb{B}) \odot D(\mathbf{a}|\mathbb{B})}{D(\mathbf{a})}, \quad \text{where } D(\mathbf{a}) \stackrel{\text{TP2}}{=} \text{tr}(D(\mathbb{B}) \odot D(\mathbf{a}|\mathbb{B})).$$

This is the Bayes rule derived in Warmuth (2005) that was discussed above.

$$\text{BR4 } D(\mathbf{b}|\mathbf{a}) = \frac{D(\mathbf{b})D(\mathbf{a}|\mathbf{b})}{D(\mathbf{a})}, \quad \text{where } D(\mathbf{a}) \stackrel{\text{TP1}}{=} \sum_i D(\mathbf{b}_i)D(\mathbf{a}|\mathbf{b}_i).$$

The summation in the normalization factor proceeds over any orthogonal system \mathbf{b}_i .

All these Bayes rules can be easily derived as follows: first express the conditional on the left i.t.o. the joint by applying the definitions of conditional probability from Sect. 7; then

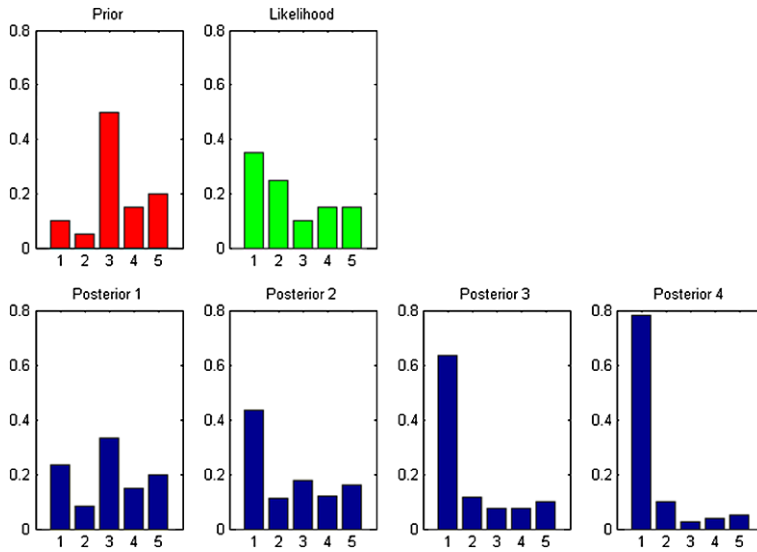


Fig. 7 (Color online) We apply the conventional Bayes rule 4 times, using the same data likelihood vector $P(y|M_i)$ and making the current posterior the new prior. At first, the posteriors are close to the initial prior but eventually the posteriors focus their weight on $\arg \max_i P(y|M_i)$. The conventional Bayes rule may be seen as a soft maximum calculation. The initial prior is in *red*, the likelihood is in *green* and posteriors are in *blue*

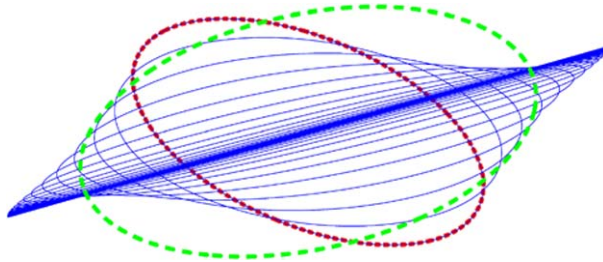


Fig. 8 (Color online) We depict several iterations of the generalized Bayes rule. The *red ellipse* depicts the prior $D(\mathbb{M})$, the *green ellipse* depicts the data likelihood matrix $D(y|\mathbb{M})$, which is kept fixed on successive iterations, and the *blue ellipses* depict posteriors $D(\mathbb{M}|y)$. The posterior density matrices gradually move away from the prior and focus on the longest axis of the covariance matrix. The generalized Bayes rule can be seen as a soft calculation of eigenvector with largest eigenvalue

apply these definitions again for expressing the joint in terms of the reverse conditional. For example,

$$D(\mathbb{B}|a) \stackrel{\text{CP2}}{=} \frac{D(\mathbb{B}, a)}{D(a)} \stackrel{\text{CP3}}{=} \frac{D(\mathbb{B}) \odot D(a|\mathbb{B})}{D(a)}.$$

As was mentioned above, the new Bayes rule can be seen as a soft maximum eigenvalue calculation. We will now give an example that shows that its impossible to track the maximum eigenvalue without changing the eigensystem. First, suppose that we have a diagonal density matrix $W = \sum_i \omega_i e_i e_i^\top$ and another diagonal matrix $S = \sum_i \sigma_i e_i e_i^\top$. Then $\text{tr}(WS) = \sum_i \omega_i \sigma_i$ and this means that by changing ω_i we can easily focus on the

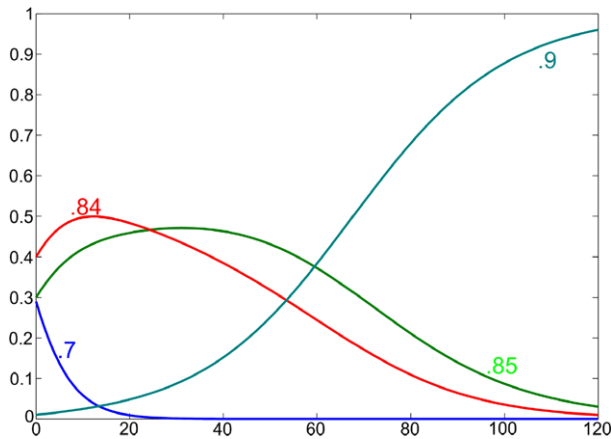


Fig. 9 We plot many iterations of the conventional Bayes rule when the same data likelihood $(P(y|M_i)) = (.7, .84, .85, .9)$ is used in each iteration and the prior is $(P(M_i)) = (.29, .4, .3, .01)$. For each of the four models we plot the posterior probability as a function of the iteration number. Initially the posterior curve with likelihood .85 overtakes the curve with likelihood .84, but eventually the curve with likelihood .9 takes over both. Note that the curve with the largest data likelihood looks like a sigmoid and the one with smallest like a reverse sigmoid

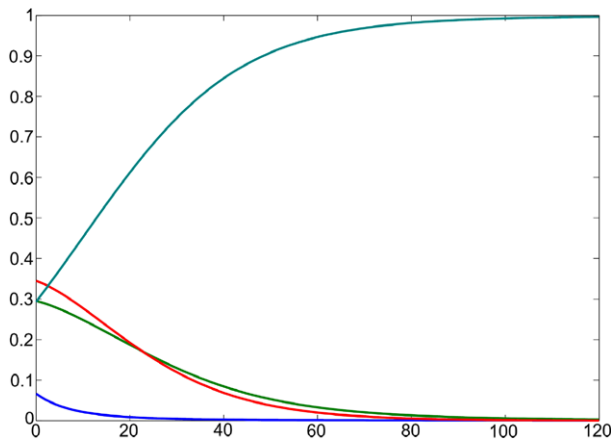


Fig. 10 We plot many iterations of the generalized Bayes rule when the same data likelihood matrix $D(y|M)$ is used in each iteration. As the prior $D(M)$ we choose the diagonalized prior $\text{diag}((P(M_i)))$ of Fig. 9 on the left and as the likelihood $D(y|M)$ we choose $U \text{diag}((P(y|M_i)))U^T$, where the eigensystem U is a random rotation matrix. Let $D(M|t)$ denote the posterior at iteration t when the fixed $D(y|M)$ is used in all iterations. The curves are the projections of this posterior onto the four eigendirections of $D(y|M)$ as a function of t , i.e. $u_i^T D(M|t) u_i$, where u_i are the columns of U . The above plot is qualitatively similar to the left plot. The curve corresponding to the largest eigenvalue of the data likelihood is again a partial sigmoid

high σ_i . Now suppose W is diagonal as before, but S has the Hadamard matrix eigensystem. Hadamard matrices H are square $n \times n$ matrices that have ± 1 elements and satisfy the condition $HH^T = nI$. Thus $\frac{H}{\sqrt{n}}$ is an orthogonal matrix. Let h_i be the columns of this orthogonal matrix derived from a Hadamard matrix and let $S = \sum_i \sigma_i h_i h_i^T$. Entries of h_i

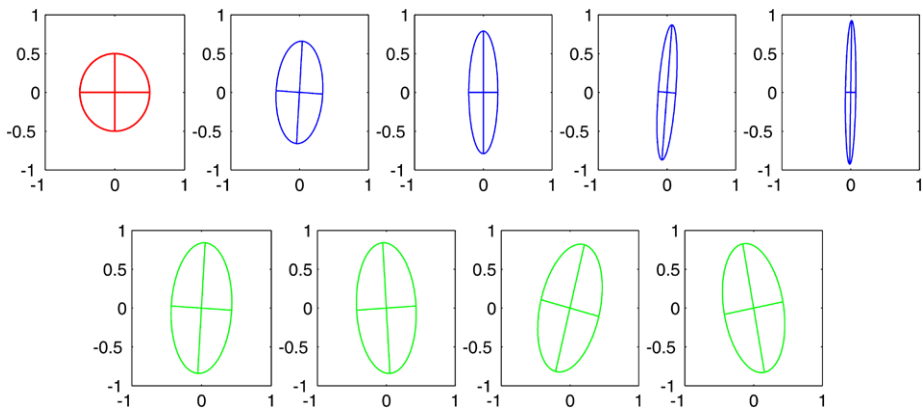


Fig. 11 (Color online) Sequence of Bayes updates with the new Bayes rule (9.2): from left to right, the prior is in red; the first data likelihood matrix is below in green; the first posterior is above in blue, and so forth

are $\pm \frac{1}{\sqrt{n}}$, therefore $\text{tr}(\mathbf{e}_i \mathbf{e}_i^\top \mathbf{h}_j \mathbf{h}_j^\top) = \frac{1}{n}$. Computing the trace we obtain:

$$\text{tr}(\mathbf{W}\mathbf{S}) = \sum_{i,j} \sigma_i \tau_j \text{tr}(\mathbf{e}_i \mathbf{e}_i^\top \mathbf{h}_j \mathbf{h}_j^\top) = \frac{1}{n} \text{tr}(\mathbf{W}) \text{tr}(\mathbf{S}) = \frac{\text{tr}(\mathbf{S})}{n}.$$

This means that *any* diagonal density matrix \mathbf{W} only “sees” the average of eigenvalues of \mathbf{S} and is unable to focus on the highest eigenvalue.

9.1 Deriving the conventional and generalized Bayes rule

In this section we show how to derive the conventional Bayes rule (9.1) and the generalized Bayes rule for density matrices (9.2) by minimizing a tradeoff between a relative entropy and an expected log likelihood. For two probability vectors \mathbf{x} and \mathbf{y} , the relative entropy is defined as $\Delta(\mathbf{x}, \mathbf{y}) := \sum_i x_i \log \frac{x_i}{y_i}$. We use the convention that $0 \log 0 := 0$ which is justified by $\lim_{x \rightarrow 0} x \log x = 0$. It is well known that $\Delta(\mathbf{x}, \mathbf{y}) \geq 0$ and that $\Delta(\mathbf{x}, \mathbf{y}) = 0$ iff $\mathbf{x} = \mathbf{y}$.

Theorem 3 *Let the prior $(P(M_i))$ be any probability vector and the data likelihood $(P(y|M_i))$ be any non-negative vector of the same dimension. Then*

$$-\log P(y) = \inf_{(\omega_i) \text{ prob. vec.}} \Delta((\omega_i), (P(M_i))) - \sum_i \omega_i \log P(y|M_i),$$

and $\boldsymbol{\omega} = (P(M_i)P(y|M_i)/P(y))$ is the unique optimum solution.

Proof Let the support of a vector \mathbf{x} be the set of all indices $1 \leq i \leq n$ s.t. $x_i \neq 0$ and denote this set as $s(\mathbf{x})$. For any probability vector (ω_i) , such that $s((\omega_i)) \subseteq s(P(M_i)) \cap s(P(y|M_i))$, we have

$$-\log P(y) = \underbrace{\sum_i \omega_i \log \frac{\omega_i}{P(M_i)}}_{\Delta((\omega_i), (P(M_i)))} - \sum_i \omega_i \log P(y|M_i) - \underbrace{\sum_i \omega_i \log \frac{\omega_i}{P(M_i)P(y|M_i)/P(y)}}_{\Delta((\omega_i), (P(M_i)P(y|M_i)/P(y)))}.$$

The precondition on the support of (ω_i) assures that all three sums above are finite because it avoids the case $\omega_i \log 0$, when $\omega_i > 0$. Since the l.h.s. is a constant,

$$\begin{aligned} & \inf_{\substack{(\omega_i) \text{ prob. vec.} \\ s((\omega_i)) \subseteq s(P(M_i)) \cap s(P(y|M_i))}} \Delta((\omega_i), (P(M_i))) - \sum_i \omega_i P(y|M_i) \\ &= \sup_{\substack{(\omega_i) \text{ prob. vec.} \\ s((\omega_i)) \subseteq s(P(M_i)) \cap s(P(y|M_i))}} -\Delta((\omega_i), (P(M_i)P(y|M_i)/P(y))). \end{aligned}$$

The sup clearly has $\omega = (P(M_i)P(y|M_i)/P(y))$ as its unique solution and the inf remains unchanged if the condition on the support of (ω_i) is dropped. This gives us the statement of the theorem. \square

This theorem can also be proven using differentiation (see e.g. Zellner 1998; Kivinen and Warmuth 1997; Singh et al. 2003). For the density matrix case this was done in Warmuth (2005), Tsuda et al. (2005). We now prove the corresponding theorem for density matrices in a different way. For two density matrices \mathbf{A} and \mathbf{B} , the quantum relative entropy is defined as $\Delta(\mathbf{A}, \mathbf{B}) := \text{tr}(\mathbf{A}(\log \mathbf{A} - \log \mathbf{B}))$. There is a potential problem when some of the eigenvalues of the matrices are zero. However, we will now reason that this definition is justified under the assumption $0 \log 0 = 0$ and $\Delta(\mathbf{A}, \mathbf{B})$ is bounded iff $\text{range}(\mathbf{A}) \subseteq \text{range}(\mathbf{B})$.

The first term $\text{tr}(\mathbf{A} \log \mathbf{A})$ becomes $\sum_i \alpha_i \log \alpha_i$, where the α_i are the eigenvalues of \mathbf{A} . This term is always finite. If \mathbf{B} is eigendecomposed as $\sum_i \beta_i \mathbf{b}_i \mathbf{b}_i^\top$, then the second term $\text{tr}(\mathbf{A} \log \mathbf{B})$ can be rewritten as $\sum_i \mathbf{b}_i^\top \mathbf{A} \mathbf{b}_i \log \beta_i$. If $\text{range}(\mathbf{A}) \subseteq \text{range}(\mathbf{B})$, then $\text{range}(\mathbf{B})^\perp \subseteq \text{range}(\mathbf{A})^\perp$, where \perp denotes the orthogonal complement space. If $\beta_i = 0$, then $\mathbf{b}_i \in \text{range}(\mathbf{B})^\perp$ and under our assumption on $\text{range}(\mathbf{A})$ this also means that $\mathbf{b}_i^\top \mathbf{A} \mathbf{b}_i = 0$. Therefore, for all i , s.t. $\beta_i = 0$, the summand $\mathbf{b}_i^\top \mathbf{A} \mathbf{b}_i \log \beta_i$ has the form $0 \log 0 = 0$. If on the other hand, $\text{range}(\mathbf{A}) \not\subseteq \text{range}(\mathbf{B})$, this also means $\text{range}(\mathbf{B})^\perp \not\subseteq \text{range}(\mathbf{A})^\perp$. The eigenvectors \mathbf{b}_i with zero eigenvalues form a basis for $\text{range}(\mathbf{B})^\perp$ and therefore there exists some \mathbf{b}_i s.t. $\mathbf{b}_i^\top \mathbf{A} \mathbf{b}_i \neq 0$. This gives a summand of the form $x \log 0$, with $x \neq 0$, and this is infinite. Notice that this discussion also means that

$$\text{tr}(\mathbf{A} \log \mathbf{B}) = \begin{cases} \text{tr}(\mathbf{A} \log^+ \mathbf{B}) & \text{when } \text{range}(\mathbf{A}) \subseteq \text{range}(\mathbf{B}), \\ -\infty & \text{otherwise.} \end{cases} \quad (9.3)$$

As before the function $\Delta(\mathbf{A}, \mathbf{B})$ is non-negative and equal zero iff both arguments agree (e.g. Nielsen and Chuang 2000).

Theorem 4 *Let the prior $\mathbf{D}(\mathbb{M})$ be any density matrix and data likelihood $\mathbf{D}(\mathbf{y}|\mathbb{M})$ be any symmetric positive definite matrix of the same dimension. Then*

$$-\log \mathbf{D}(\mathbf{y}) = \inf_{\mathbf{W} \text{ dens. mat.}} \Delta(\mathbf{W}, \mathbf{D}(\mathbb{M})) - \text{tr}(\mathbf{W} \log \mathbf{D}(\mathbf{y}|\mathbb{M})),$$

and $\mathbf{W} = \frac{\mathbf{D}(\mathbb{M}) \odot \mathbf{D}(\mathbf{y}|\mathbb{M})}{\mathbf{D}(\mathbf{y})}$ is the unique optimum solution.

Proof For any density matrix \mathbf{W} s.t. $\text{range}(\mathbf{W}) \subseteq \text{range}(\mathbf{D}(\mathbb{M})) \cap \text{range}(\mathbf{D}(\mathbf{y}|\mathbb{M}))$, we have

$$\begin{aligned} -\log \mathbf{D}(\mathbf{y}) &= \underbrace{\text{tr}(\mathbf{W}(\log \mathbf{W} - \log \mathbf{D}(\mathbb{M})))}_{\Delta(\mathbf{W}, \mathbf{D}(\mathbb{M}))} - \text{tr}(\mathbf{W}(\log \mathbf{D}(\mathbf{y}|\mathbb{M}))) \\ &\quad - \text{tr}(\mathbf{W}(\log \mathbf{W} - (\log \mathbf{D}(\mathbb{M}) + \log \mathbf{D}(\mathbf{y}|\mathbb{M})/\mathbf{D}(\mathbf{y}))))). \end{aligned}$$

Since $\text{range}(\mathbf{W}) \subseteq \text{range}(\mathbf{D}(\mathbb{M})) \cap \text{range}(\mathbf{D}(\mathbf{y}|\mathbb{M}))$, $\text{tr}(\mathbf{W} \log \mathbf{D}(\mathbb{M}))$ and $\text{tr}(\mathbf{W} \log \mathbf{D}(\mathbf{y}|\mathbb{M}))$ are both finite. Assuming that for any symmetric positive definite matrices \mathbf{W} , \mathbf{A} and \mathbf{B}

$$\text{tr}(\mathbf{W}(\log \mathbf{A} + \log \mathbf{B})) = \text{tr}(\mathbf{W} \log(\mathbf{A} \odot \mathbf{B})), \text{ when } \text{range}(\mathbf{W}) \subseteq \text{range}(\mathbf{A}) \cap \text{range}(\mathbf{B}), \quad (9.4)$$

the above equality would become

$$-\log D(\mathbf{y}) = \Delta(\mathbf{W}, \mathbf{D}(\mathbb{M})) - \text{tr}(\mathbf{W} \log \mathbf{D}(\mathbf{y}|\mathbb{M})) - \Delta(\mathbf{W}, (\mathbf{D}(\mathbb{M}) \odot \mathbf{D}(\mathbf{y}|\mathbb{M})) / \mathbf{D}(\mathbf{y})).$$

Since the l.h.s. is a constant,

$$\begin{aligned} & \inf_{\substack{\mathbf{W} \text{ dens.mat.} \\ \text{range}(\mathbf{W}) \subseteq \text{range}(\mathbf{D}(\mathbb{M})) \cap \text{range}(\mathbf{D}(\mathbf{y}|\mathbb{M}))}} \Delta(\mathbf{W}, \mathbf{D}(\mathbb{M})) - \text{tr}(\mathbf{W} \log \mathbf{D}(\mathbf{y}|\mathbb{M})) \\ &= \sup_{\substack{\mathbf{W} \text{ dens.mat.} \\ \text{range}(\mathbf{W}) \subseteq \text{range}(\mathbf{D}(\mathbb{M})) \cap \text{range}(\mathbf{D}(\mathbf{y}|\mathbb{M}))}} -\Delta(\mathbf{W}, (\mathbf{D}(\mathbb{M}) \odot \mathbf{D}(\mathbf{y}|\mathbb{M})) / \mathbf{D}(\mathbf{y})). \end{aligned}$$

The sup clearly has the unique solution $\frac{\mathbf{D}(\mathbb{M}) \odot \mathbf{D}(\mathbf{y}|\mathbb{M})}{\mathbf{D}(\mathbf{y})}$ and the inf remains unchanged if the condition on the range of \mathbf{W} is dropped. This gives us the statement of the theorem.

We still need to show (9.4). Since $\text{range}(\mathbf{W}) \subseteq \text{range}(\mathbf{A}) \cap \text{range}(\mathbf{B}) \stackrel{\text{OP1}}{=} \text{range}(\mathbf{A} \odot \mathbf{B})$,

$$\begin{aligned} \text{tr}(\mathbf{W} \log(\mathbf{A} \odot \mathbf{B})) &\stackrel{(9.3)}{=} \text{tr}(\mathbf{W} \log^+(\mathbf{A} \odot \mathbf{B})) \\ &\stackrel{(4.6)}{=} \text{tr}(\mathbf{W} \mathbf{P}_{A \cap B} (\log^+ \mathbf{A} + \log^+ \mathbf{B}) \mathbf{P}_{A \cap B}) \\ &= \text{tr}(\underbrace{\mathbf{P}_{A \cap B} \mathbf{W} \mathbf{P}_{A \cap B}}_{\mathbf{W}} (\log^+ \mathbf{A} + \log^+ \mathbf{B})) \\ &\stackrel{(9.3)}{=} \text{tr}(\mathbf{W} (\log \mathbf{A} + \log \mathbf{B})). \quad \square \end{aligned}$$

We conclude with a discussion of the relationship between the conventional Bayes rule for probability vectors and the generalized Bayes rule for density matrices. Density matrices are determined by a probability vector of eigenvalues as well as an orthogonal eigensystem. An orthogonal system \mathbf{w}_i turns the prior density matrix $\mathbf{D}(\mathbb{M})$ into the probability vector $(\text{tr}(\mathbf{D}(\mathbb{M}) \mathbf{w}_i \mathbf{w}_i^\top))$, which we call a *pinching* of $\mathbf{D}(\mathbb{M})$. Similarly the pinching of the data likelihood matrix $\mathbf{D}(\mathbf{y}|\mathbb{M})$ is the vector $(\text{tr}(\mathbf{D}(\mathbf{y}|\mathbb{M}) \mathbf{w}_i \mathbf{w}_i^\top)) \in [0, 1]^n$. The idea is to express our Bayes rule for density matrices as the conventional Bayes rule for the pinched priors and likelihoods w.r.t. a certain eigensystem. That is, we want to be able to say that the generalized Bayes rule is the conventional Bayes rule for the “best” pinching.

The above outline is essentially true, but we need to pinch in the log domain. With (9.3), Property OP11 can be extended to

$$\text{tr}(\mathbf{A} \odot \mathbf{u} \mathbf{u}^\top) = e^{\mathbf{u}^\top \log \mathbf{A} \mathbf{u}}, \quad \text{for any unit } \mathbf{u} \text{ and symmetric positive definite matrix } \mathbf{A}. \quad (9.5)$$

We call $\text{tr}(\mathbf{A} \odot \mathbf{w}_i \mathbf{w}_i^\top)$ a *remote pinching* of \mathbf{A} . Since its components satisfy $\text{tr}(\mathbf{D}(\mathbb{M}) \odot \mathbf{w}_i \mathbf{w}_i^\top) \stackrel{\text{OP10}}{\leq} \text{tr}(\mathbf{D}(\mathbb{M}) \mathbf{w}_i \mathbf{w}_i^\top)$, the remote pinchings of $\mathbf{D}(\mathbb{M})$ must be normalized to form a probability vector.

We can rewrite the argument of the optimization problem for the generalized Bayes rule based on the eigendecomposition $\mathcal{W} \mathcal{W}^\top$ of the density matrix \mathbf{W} :

$$\begin{aligned} & \Delta(\mathbf{W}, \mathbf{D}(\mathbb{M})) - \text{tr}(\mathbf{W} \log \mathbf{D}(\mathbf{y}|\mathbb{M})) \\ &= \text{tr}(\omega \mathcal{W}^\top (\log \mathbf{W} - \log \mathbf{D}(\mathbb{M})) \mathcal{W}) - \text{tr}(\omega \mathcal{W}^\top (\log \mathbf{D}(\mathbf{y}|\mathbb{M})) \mathcal{W}) \end{aligned}$$

$$\begin{aligned}
&= \sum_i \omega_i (\log \omega_i - \mathbf{w}_i^\top (\log \mathbf{D}(\mathbb{M})) \mathbf{w}_i) - \sum_i \omega_i \mathbf{w}_i^\top (\log \mathbf{D}(\mathbf{y}|\mathbb{M})) \mathbf{w}_i \\
&\stackrel{(9.5)}{=} \sum_i \omega_i (\log \omega_i - \log \text{tr}(\mathbf{D}(\mathbb{M}) \odot \mathbf{w}_i \mathbf{w}_i^\top)) - \sum_i \omega_i \log \text{tr}(\mathbf{D}(\mathbf{y}|\mathbb{M}) \odot \mathbf{w}_i \mathbf{w}_i^\top) \\
&= \Delta((\omega_i), (\text{tr}(\mathbf{D}(\mathbb{M}) \odot \mathbf{w}_i \mathbf{w}_i^\top) / Z_{\mathcal{W}})) - \sum_i \omega_i \log \text{tr}(\mathbf{D}(\mathbf{y}|\mathbb{M}) \odot \mathbf{w}_i \mathbf{w}_i^\top) - \log Z_{\mathcal{W}},
\end{aligned}$$

where the normalization $Z_{\mathcal{W}} = \sum_j \text{tr}(\mathbf{D}(\mathbb{M}) \odot \mathbf{w}_j \mathbf{w}_j^\top)$ does not depend on the eigenvalues. By Theorem 3, the above is minimized w.r.t. ω when $\omega = (P_{\mathcal{W}}(M_i) P_{\mathcal{W}}(\mathbf{y}|M_i) / P_{\mathcal{W}}(\mathbf{y}))$, where $P_{\mathcal{W}}(M_i) := \text{tr}(\mathbf{D}(\mathbb{M}) \odot \mathbf{w}_i \mathbf{w}_i^\top) / Z_{\mathcal{W}}$ is the normalized remote pinching of the prior and $P_{\mathcal{W}}(\mathbf{y}|M_i) := \text{tr}(\mathbf{D}(\mathbf{y}|\mathbb{M}) \odot \mathbf{w}_i \mathbf{w}_i^\top)$ is the remote pinching of the data likelihood matrix. With this optimum choice of ω , the minimization problem of the generalized Bayes rule simplifies to

$$\begin{aligned}
&\inf_{\mathcal{W} \mathcal{W}^\top = I} -\log P_{\mathcal{W}}(\mathbf{y}) - \log Z_{\mathcal{W}} \\
&= \inf_{\mathcal{W} \mathcal{W}^\top = I} -\log \left(\sum_i \text{tr}(\mathbf{D}(\mathbb{M}) \odot \mathbf{w}_i \mathbf{w}_i^\top) \text{tr}(\mathbf{D}(\mathbf{y}|\mathbb{M}) \odot \mathbf{w}_i \mathbf{w}_i^\top) \right) \\
&\stackrel{\text{OP15}}{=} \inf_{\mathcal{W} \mathcal{W}^\top = I} -\log \left(\sum_i \text{tr}((\mathbf{D}(\mathbb{M}) \odot \mathbf{D}(\mathbf{y}|\mathbb{M})) \odot \mathbf{w}_i \mathbf{w}_i^\top) \right) \\
&\stackrel{\text{OP10}}{\geq} \inf_{\mathcal{W} \mathcal{W}^\top = I} -\log \left(\sum_i \text{tr}((\mathbf{D}(\mathbb{M}) \odot \mathbf{D}(\mathbf{y}|\mathbb{M})) \mathbf{w}_i \mathbf{w}_i^\top) \right) \\
&= -\log \text{tr}(\mathbf{D}(\mathbb{M}) \odot \mathbf{D}(\mathbf{y}|\mathbb{M})).
\end{aligned}$$

The above inequality is tight iff \mathcal{W} is an eigensystem of $\mathbf{D}(\mathbb{M}) \odot \mathbf{D}(\mathbf{y}|\mathbb{M})$. We conclude that the optimization problem for the generalized Bayes rule is optimized when \mathcal{W} is an eigensystem of $\mathbf{D}(\mathbb{M}) \odot \mathbf{D}(\mathbf{y}|\mathbb{M})$ and the vector of eigenvalues ω is conventional posterior derived from the normalized remote pinchings of the prior and the remote pinchings of the data likelihood.

9.2 Chaining of the Bayes rule

The conventional Bayes rule can be applied iteratively to a sequence of data and various cancellations occur. For the sake of simplicity we only consider two data points y_1, y_2 :

$$P(M_i | y_2, y_1) = \frac{P(M_i | y_1) P(y_2 | M_i, y_1)}{P(y_2 | y_1)} = \frac{P(M_i) P(y_1 | M_i) P(y_2 | M_i, y_1)}{P(y_2 | y_1) P(y_1)}.$$

The normalization can be rewritten as:

$$\begin{aligned} P(y_2|y_1)P(y_1) &= \left(\sum_i \underbrace{P(M_i|y_1)}_{\text{using (9.1)}} P(y_2|M_i, y_1) \right) \left(\sum_i P(M_i) P(y_1|M_i) \right) \\ &= \sum_i P(M_i) P(y_1|M_i) P(y_2|M_i, y_1) = P(y_2, y_1). \end{aligned} \quad (9.6)$$

Analogously, by essentially applying the generalized Bayes rule (9.2) two times we get:

$$D(\mathbb{M}|y_2, y_1) = \frac{D(\mathbb{M}|y_1) \odot D(y_2|\mathbb{M}, y_1)}{D(y_2|y_1)} = \frac{D(\mathbb{M}) \odot D(y_1|\mathbb{M}) \odot D(y_2|\mathbb{M}, y_1)}{D(y_2|y_1)D(y_1)}.$$

As in the diagonal case (9.6), the normalization can be rewritten into one term (by applying TP2 twice and then the generalized Bayes rule (9.2)):

$$\begin{aligned} D(y_2|y_1)D(y_1) &= \text{tr}(D(\mathbb{M}|y_1) \odot D(y_2|\mathbb{M}, y_1)) \text{tr}(D(\mathbb{M}) \odot D(y_1|\mathbb{M})) \\ &= \text{tr} \left(\frac{D(\mathbb{M}) \odot D(y_1|\mathbb{M})}{\text{tr}(D(\mathbb{M}) \odot D(y_1|\mathbb{M}))} \odot D(y_2|\mathbb{M}, y_1) \right) \text{tr}(D(\mathbb{M}) \odot D(y_1|\mathbb{M})) \\ &= \text{tr}(D(\mathbb{M}) \odot D(y_1|\mathbb{M}) \odot D(y_2|\mathbb{M}, y_1)) = D(y_1, y_2). \end{aligned}$$

Finally as in (8.1), we can upper bound the data probability $D(y_1, y_2)$ in terms of the product of the expected variances for the two trials:

$$\begin{aligned} D(y_2, y_1) &= \text{tr}(D(\mathbb{M}|y_1) \odot D(y_2|\mathbb{M}, y_1)) \text{tr}(D(\mathbb{M}) \odot D(y_1|\mathbb{M})) \\ &\leq \text{tr}(D(\mathbb{M}|y_1)D(y_2|\mathbb{M}, y_1)) \text{tr}(D(\mathbb{M})D(y_1|\mathbb{M})). \end{aligned}$$

9.3 Bounds

Recall the following conventional bound for the negative log-likelihood of the data i.t.o. the negative log-likelihood of the MAP estimator:

$$\begin{aligned} -\log P(y) &= -\log \sum_i P(y|M_i)P(M_i) \\ &\leq \min_i (-\log P(y|M_i) - \log P(M_i)). \end{aligned} \quad (9.7)$$

We will give analogous bound for density matrices. For this we need the following inequality: For any unit vector \mathbf{m} and symmetric positive definite matrix \mathbf{A} :

$$-\log \mathbf{m}^\top \mathbf{A} \mathbf{m} \stackrel{\text{OP10}}{\leq} -\log \text{tr}(\mathbf{A} \odot \mathbf{m} \mathbf{m}^\top) \stackrel{(9.5)}{=} -\mathbf{m}^\top (\log \mathbf{A}) \mathbf{m}. \quad (9.8)$$

Using the fact that $\text{tr}(\mathbf{A}) \geq \mathbf{m}^\top \mathbf{A} \mathbf{m}$, we can now prove an analogous MAP bound for the generalized probabilities:

$$\begin{aligned} -\log D(\mathbf{y}) &= -\log \text{tr}(D(\mathbf{y}|\mathbb{M}) \odot D(\mathbb{M})) \\ &\leq \min_{\mathbf{m}} (-\log \mathbf{m}^\top (D(\mathbf{y}|\mathbb{M}) \odot D(\mathbb{M})) \mathbf{m}) \end{aligned}$$

$$\begin{aligned}
&\stackrel{(9.8)}{\leq} \min_m (-\mathbf{m}^\top \log(\mathbf{D}(\mathbf{y}|\mathbb{M}) \odot \mathbf{D}(\mathbb{M}))\mathbf{m}) \\
&\leq \min_m (-\mathbf{m}^\top \log \mathbf{D}(\mathbf{y}|\mathbb{M})\mathbf{m} - \mathbf{m}^\top \log \mathbf{D}(\mathbb{M})\mathbf{m}).
\end{aligned}$$

The last inequality becomes (9.4), when $\mathbf{m} \in \text{range}(\mathbf{D}(\mathbb{M})) \cap \text{range}(\mathbf{D}(\mathbf{y}|\mathbb{M}))$. Otherwise, it holds trivially because $-\mathbf{m}^\top \log(\mathbf{D}(\mathbf{y}|\mathbb{M}) \odot \mathbf{D}(\mathbb{M}))\mathbf{m} = +\infty$.

Intuitively, there are two domains: the probability domain and the log probability domain. The conventional bound (9.7) can also be written in the probability domain:

$$P(\mathbf{y}) \geq \max_i P(M_i)P(\mathbf{y}|M_i).$$

However for the generalized probability case, there does not seem to be a simple similar inequality in the probability domain. Throughout the paper we always notice that the matrix operations need to be done in the log domain.

In the conventional case $P(\mathbf{y})$ is also upper bounded by $\max_i P(\mathbf{y}|M_i)$. For the generalized case, the analogous formula is the following, where μ_i and \mathbf{m}_i are the eigenvalues/vectors of $\mathbf{D}(\mathbb{M})$ and \mathbf{m} any unit direction:

$$\begin{aligned}
\mathbf{D}(\mathbf{y}) &= \text{tr}(\mathbf{D}(\mathbf{y}|\mathbb{M}) \odot \mathbf{D}(\mathbb{M})) \\
&\leq \text{tr}(\mathbf{D}(\mathbf{y}|\mathbb{M})\mathbf{D}(\mathbb{M})) \\
&= \sum_i \mu_i \mathbf{m}_i^\top \mathbf{D}(\mathbf{y}|\mathbb{M})\mathbf{m}_i \\
&\leq \max_i \mathbf{m}_i^\top \mathbf{D}(\mathbf{y}|\mathbb{M})\mathbf{m}_i \\
&\leq \max_m \mathbf{m}^\top \mathbf{D}(\mathbf{y}|\mathbb{M})\mathbf{m}.
\end{aligned}$$

10 Summary of the probability calculus for density matrices

In this section we give a summary of all the rules of our calculus. The definitions are indicated with $:=$ and at the end we summarize the justification for our choice of definitions. Table 1 shows connections between different objects and the formulas that relate them.

10.1 Marginalization rules for joints of Sects. 5 and 6

$$\begin{aligned}
\text{MJ1 } \mathbf{D}(\mathbf{a}) &:= \text{tr}(\mathbf{D}(\mathbb{A})\mathbf{a}\mathbf{a}^\top) = \mathbf{a}^\top \mathbf{D}(\mathbb{A})\mathbf{a}. \\
\text{MJ2 } \mathbf{D}(\mathbb{A}) &:= \text{tr}_{\mathbb{B}}(\mathbf{D}(\mathbb{A}, \mathbb{B})). \\
\text{MJ3 } \mathbf{D}(\mathbf{a}, \mathbf{b}) &:= \text{tr}(\mathbf{D}(\mathbb{A}, \mathbb{B})(\mathbf{a} \otimes \mathbf{b})(\mathbf{a} \otimes \mathbf{b})^\top) = \text{tr}(\mathbf{D}(\mathbb{A}, \mathbb{B})(\mathbf{a}\mathbf{a}^\top \otimes \mathbf{b}\mathbf{b}^\top)). \\
\text{MJ4 } \mathbf{D}(\mathbb{A}, \mathbf{b}) &:= \text{tr}_{\mathbb{B}}(\mathbf{D}(\mathbb{A}, \mathbb{B})(\mathbf{I}_{\mathbb{A}} \otimes \mathbf{b}\mathbf{b}^\top)). \\
\text{MJ5 } \mathbf{D}(\mathbf{a}, \mathbf{b}) &= \text{tr}(\mathbf{D}(\mathbb{A}, \mathbf{b})\mathbf{a}\mathbf{a}^\top).
\end{aligned}$$

10.2 Conditional probability rules of Sect. 7

$$\begin{aligned}
\text{CP1 } \mathbf{D}(\mathbb{A}|\mathbb{B}) &:= \mathbf{D}(\mathbb{A}, \mathbb{B}) \odot (\mathbf{I}_{\mathbb{A}} \otimes \mathbf{D}(\mathbb{B}))^{-1}. \\
\text{CP2 } \mathbf{D}(\mathbb{A}|\mathbf{b}) &:= \frac{\mathbf{D}(\mathbb{A}, \mathbf{b})}{\text{tr}(\mathbf{D}(\mathbb{A}, \mathbf{b}))}. \\
\text{CP3 } \mathbf{D}(\mathbf{a}|\mathbb{B}) &:= \mathbf{D}(\mathbf{a}, \mathbb{B}) \odot \mathbf{D}(\mathbb{B})^{-1}. \\
\text{CP4 } \mathbf{D}(\mathbf{a}|\mathbf{b}) &:= \frac{\mathbf{D}(\mathbf{a}, \mathbf{b})}{\mathbf{D}(\mathbf{b})}.
\end{aligned}$$

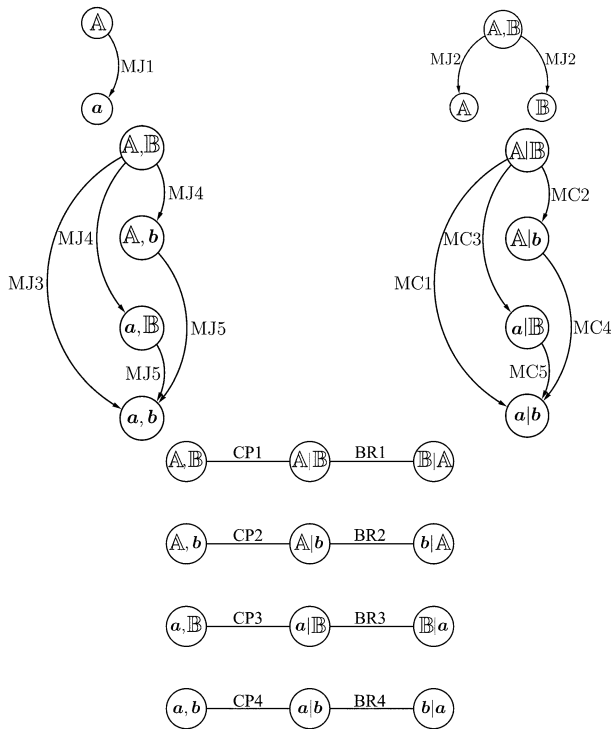


Table 1 A series of charts summarizing the different relationships for joints and conditionals. Each edge references the formula stating the relationship. For symmetric cases only one formula is given and the corresponding edges in the chart will have the same label

CP1 has the form: density matrix \odot inverse of a normalization. Below we reexpress the other definitions in this unified form:

$$\text{CP'2 } D(\mathbb{A}|\mathbb{b}) = \text{tr}_{\mathbb{B}}(D(\mathbb{A}, \mathbb{B})(I_{\mathbb{A}} \otimes \mathbb{b}\mathbb{b}^{\top})) \odot \text{tr}_{\mathbb{B}}((I_{\mathbb{A}} \otimes D(\mathbb{B}))(I_{\mathbb{A}} \otimes \mathbb{b}\mathbb{b}^{\top}))^{-1}.$$

$$\text{CP'3 } D(\mathbb{a}|\mathbb{B}) = \text{tr}_{\mathbb{A}}(D(\mathbb{A}, \mathbb{B})(\mathbb{a}\mathbb{a}^{\top} \otimes I_{\mathbb{B}})) \odot \text{tr}_{\mathbb{A}}((I_{\mathbb{A}} \otimes D(\mathbb{B}))(\mathbb{a}\mathbb{a}^{\top} \otimes I_{\mathbb{B}}))^{-1}.$$

$$\text{CP'4 } D(\mathbb{a}|\mathbb{b}) = \text{tr}(D(\mathbb{A}, \mathbb{B})(\mathbb{a}\mathbb{a}^{\top} \otimes \mathbb{b}\mathbb{b}^{\top})) \odot \text{tr}((I_{\mathbb{A}} \otimes D(\mathbb{B}))(\mathbb{a}\mathbb{a}^{\top} \otimes \mathbb{b}\mathbb{b}^{\top}))^{-1}.$$

10.3 Marginalization rules for conditionals of Sect. 7

$$\text{MC1 } D(\mathbb{a}|\mathbb{b}) = \frac{\text{tr}(D(\mathbb{A}|\mathbb{B}) \odot (I_{\mathbb{A}} \otimes D(\mathbb{B}))(\mathbb{a}\mathbb{a}^{\top} \otimes \mathbb{b}\mathbb{b}^{\top}))}{\text{tr}(D(\mathbb{B})\mathbb{b}\mathbb{b}^{\top})}.$$

$$\text{MC2 } D(\mathbb{A}|\mathbb{b}) = \frac{\text{tr}_{\mathbb{B}}(D(\mathbb{A}|\mathbb{B}) \odot (I_{\mathbb{A}} \otimes D(\mathbb{B}))(\mathbb{A}\mathbb{A}^{\top} \otimes \mathbb{b}\mathbb{b}^{\top}))}{\text{tr}(D(\mathbb{B})\mathbb{b}\mathbb{b}^{\top})}.$$

$$\text{MC3 } D(\mathbb{a}|\mathbb{B}) = \text{tr}_{\mathbb{A}}(D(\mathbb{A}|\mathbb{B}) \odot (I_{\mathbb{A}} \otimes D(\mathbb{B}))(\mathbb{a}\mathbb{a}^{\top} \otimes I_{\mathbb{B}})) \odot D(\mathbb{B})^{-1}.$$

$$\text{MC4 } D(\mathbb{a}|\mathbb{b}) = \text{tr}(D(\mathbb{A}|\mathbb{b})\mathbb{a}\mathbb{a}^{\top}).$$

$$\text{MC5 } D(\mathbb{a}|\mathbb{b}) = \frac{\text{tr}((D(\mathbb{a}|\mathbb{B}) \odot D(\mathbb{B}))\mathbb{b}\mathbb{b}^{\top})}{\text{tr}(D(\mathbb{B})\mathbb{b}\mathbb{b}^{\top})}.$$

All the rules here except for MC4 require additional information for marginalization, which was not necessary in the conventional case. See discussion of marginalization of conditionals in Sect. 7.

10.4 Theorems of total probability of Sect. 8

TP1 $D(\mathbf{a}) = \sum_i D(\mathbf{a}|\mathbf{b}_i)D(\mathbf{b}_i)$ for any orthogonal system \mathbf{b}_i of space \mathbb{B} .

TP2 $D(\mathbb{A}) = \sum_i D(\mathbb{A}, \mathbf{b}_i)$ for any orthogonal system \mathbf{b}_i of space \mathbb{B} .

TP3 $D(\mathbb{A}) = \text{tr}_{\mathbb{B}}(D(\mathbb{A}|\mathbb{B}) \odot (I_{\mathbb{A}} \otimes D(\mathbb{B})))$.

10.5 Bayes rules of Sect. 9

BR1 $D(\mathbb{B}|\mathbb{A}) = (I_{\mathbb{A}} \otimes D(\mathbb{B})) \odot D(\mathbb{A}|\mathbb{B}) \odot (D(\mathbb{A}) \otimes I_{\mathbb{B}})^{-1}$, where $D(\mathbb{A}) = \text{tr}_{\mathbb{B}}((I_{\mathbb{A}} \otimes D(\mathbb{B})) \odot D(\mathbb{A}|\mathbb{B}))$.

BR2 $D(\mathbf{a}|\mathbb{B}) = D(\mathbf{a})D(\mathbb{B}|\mathbf{a}) \odot D(\mathbb{B})^{-1}$, where $D(\mathbb{B}) = \text{tr}_{\mathbb{A}}(D(\mathbb{B}|\mathbb{A}) \odot (D(\mathbb{A}) \otimes I_{\mathbb{B}}))$.

BR3 $D(\mathbb{B}|\mathbf{a}) = \frac{D(\mathbb{B}) \odot D(\mathbf{a}|\mathbb{B})}{D(\mathbf{a})}$, where $D(\mathbf{a}) = \text{tr}(D(\mathbb{B}) \odot D(\mathbf{a}|\mathbb{B}))$.

BR4 $D(\mathbf{b}|\mathbf{a}) = \frac{D(\mathbf{a}|\mathbf{b})D(\mathbf{b})}{D(\mathbf{a})}$, where $D(\mathbf{a}) = \sum_i D(\mathbf{a}|\mathbf{b}_i)D(\mathbf{b}_i)$ and the summation is over any orthogonal system \mathbf{b}_i .

10.6 Summary of justifications for the definitions

Note that only the rules MJ1–4 and CP rules are definitions. Everything else in our calculus can be derived from these. MJ1 is justified by Gleason’s Theorem as discussed in Sect. 3. Gleason’s Theorem also justifies MJ3, where the Kronecker product provides the natural way to specify a joint unit (see discussion in Sect. 5). MJ2 is standard in quantum physics and $\text{tr}_{\mathbb{B}}(D(\mathbb{A}, \mathbb{B}))$ was shown to be a density matrix in Lemma 1. The rule is also compatible with the conventional case as well as with the natural generalization of independence discussed in Sect. 6. MJ4 is the natural definition of $D(\mathbb{A}, \mathbf{b})$ that satisfies MJ5 and is compatible with the conventional case.

We will outline how CP2 can be motivated as a quantum relative entropy projection. For positive definite matrices \mathbf{A} and \mathbf{B} , we extend the definition of quantum relative entropy as follows: $\Delta(\mathbf{A}, \mathbf{B}) = \text{tr}(\mathbf{A}(\log \mathbf{A} - \log \mathbf{B}) + \mathbf{B} - \mathbf{A})$. Note that this “unnormalized” relative entropy, coincides with the standard one when \mathbf{A} and \mathbf{B} have trace one. Now CP2 is motivated as

$$D(\mathbb{A}|\mathbf{b}) = \arg \inf_{W \text{ dens. mat.}} \Delta(W, D(\mathbb{A}, \mathbf{b})).$$

CP3 is motivated analogous to the generalized Bayes rule (see Sect. 9):

$$D(\mathbf{a}, \mathbf{B}) = \arg \inf_W \Delta(W, D(\mathbb{B})) - \text{tr}(W \log D(\mathbf{a}|\mathbb{B})).$$

CP1 can be motivated in a similar fashion, but now the variable is over the joint space (\mathbb{A}, \mathbb{B}) :

$$D(\mathbf{A}, \mathbf{B}) = \arg \inf_{W \text{ dens.mat.}} \Delta(W, I_{\mathbb{A}} \otimes D(\mathbb{B})) - \text{tr}(W \log D(\mathbb{A}|\mathbb{B})).$$

CP1 also was previously used in Cerf and Adami (1999) to allow a suitable definition of conditional quantum entropy. Finally, the last rule CP4 was chosen in analogy to the conventional case. It also has an interpretation as two successive quantum measurements (see Appendix).

Historically, we first justified the generalized Bayes rule BR3 based on the minimum relative entropy principle (see Warmuth 2005 and Sect. 9). After that we chose definitions CP1–CP4 to be compatible with this generalized Bayes rule.

11 Conclusions

Density matrices are central to quantum physics. We utilize many mathematical techniques from that field to develop a Bayesian probability calculus for density matrices. Intuitively, the new calculus will be useful when the data likelihood $\mathbf{D}(\mathbf{y}|\mathbb{M})$ has non-zero off-diagonal elements, i.e. information about which components are correlated or anti-correlated. The main new operation $\mathbf{A} \odot \mathbf{B}$ first takes logs of the matrices adds the logs and finally exponentiates. Any straightforward implementation of the \odot operation requires the eigendecompositions of the matrices, which are expensive to obtain. Throughout our work we notice that the log domain seems to be more important in the matrix case.

Interestingly enough the \odot operation has also been employed in computer graphics for combining affine transformation (Alexa 2002). Also the simulation of quantum computations based on the Lie Trotter Formula (Nielsen and Chuang 2000, Chap. 4.7) can be interpreted as applying the \odot operation to unitary matrices and not to symmetric positive definite matrices as we do in this paper.

The main update in quantum physics is a unitary evolution of the current density matrix \mathbf{A} , i.e. $\mathbf{A} := \mathbf{U}\mathbf{A}\mathbf{U}^\top$, where \mathbf{U} is unitary. For example, the main differential equation for density matrices in quantum physics is the following version of the Schrödinger Equation (Feynman 1972):

$$\frac{\partial \mathbf{D}(\mathbb{M}|t)}{\partial t} = i(\mathbf{H}\mathbf{D}(\mathbb{M}|t) - \mathbf{D}(\mathbb{M}|t)\mathbf{H}), \quad \text{where } \mathbf{H} \text{ is skew Hermitian.}$$

The solution has the form

$$\mathbf{D}(\mathbb{M}|t) = \exp(-it\mathbf{H})\mathbf{D}(\mathbb{M}|0)\exp(it\mathbf{H}),$$

where $\mathbf{D}(\mathbb{M}|0)$ is the initial density matrix. Since $it\mathbf{H}$ is skew Hermitian, both exponentials are unitary. Thus the above update represents a unitary transformation of the initial density matrix $\mathbf{D}(\mathbb{M}|0)$. Such transformations leave the eigenvalues unchanged and only affect the eigensystem. In contrast our generalized Bayes rule updates both the eigenvalues and eigenvectors, and the conventional Bayes rule can be seen as only updating the eigenvalues while keeping the eigenvectors fixed. Therefore the Bayes rules are decidedly not unitary updates.

For the sake of completeness we now express the Bayes rules also as solutions to differential equations. In the conventional case, the differential equations are ($1 \leq i \leq n$):

$$\frac{\partial \log P(M_i|t)}{\partial t} = \log P(y|M_i) - \sum_j P(M_j|t) \log P(y|M_j).$$

The solution is

$$P(M_i|t) = \frac{P(M_i|0)P(y|M_i)^t}{\sum_j P(M_j|0)P(y|M_j)^t}.$$

If we take the value $P(M_i|0)$ as the prior $P(M_i)$ then the expression for $P(M_i|1)$ becomes the conventional Bayes rule (9.1). There is a similar differential equation for the generalized Bayes rule (for the sake of simplicity we assume that the prior $\mathbf{D}(\mathbb{M})$ and data likelihood matrix $\mathbf{D}(\mathbf{y}|\mathbb{M})$ are strictly positive definite):

$$\frac{\partial \log \mathbf{D}(\mathbb{M}|t)}{\partial t} = \log \mathbf{D}(\mathbf{y}|\mathbb{M}) - \text{tr}(\mathbf{D}(\mathbb{M}|t) \log \mathbf{D}(\mathbf{y}|\mathbb{M})).$$

The solution has the form

$$\begin{aligned} D(\mathbb{M}|t) &= \frac{\exp(\log D(\mathbb{M}|0) + t \log D(\mathbf{y}|\mathbb{M}))}{\text{tr}(\exp(\log D(\mathbb{M}|0) + t \log D(\mathbf{y}|\mathbb{M})))} \\ &\stackrel{(4.1)}{=} \frac{D(\mathbb{M}|0) \odot D(\mathbf{y}|\mathbb{M})^t}{\text{tr}(D(\mathbb{M}|0) \odot D(\mathbf{y}|\mathbb{M})^t)}. \end{aligned}$$

If we set $D(\mathbb{M}|0)$ to the prior $D(\mathbb{M})$, then the expression for $D(\mathbb{M}|1)$ becomes the generalized Bayes rule (9.2). Notice again that the differential equations emphasize the log domain and that the \odot operation appears in the solution.

At this point we have no convincing application for the new probability calculus. However, a similar methodology was used to derive and prove bounds for parameter updates of density matrices that led to a version of Boosting (Tsuda et al. 2005) where the distribution over the examples is replaced by a density matrix, an online variance minimization algorithm where the parameter space is the unit ball (Warmuth and Kuzmin 2006), and an on-line algorithm for Principal Component Analysis (Warmuth and Kuzmin 2008).

In this paper our parameters expressing the uncertainty are symmetric positive definite matrices. However using essentially the EG_{\pm} transformation (Kivinen and Warmuth 1997), it has been shown recently that inference can be done with arbitrarily shaped matrices (Warmuth 2007). This leaves the strong possibility that the calculus developed here will generalize to arbitrary shaped matrices as well. In that case the elementary events are “asymmetric dyads” uv^T and the underlying decomposition is the SVD decomposition.

The new calculus seems to be rich enough to bring out some of the interesting phenomena of quantum physics, such as superposition and entanglement. Maybe the new calculus can be used to maintain “uncertainty” in quantum computation.

On a more technical note, we conjecture that for all non-decoupled joints $D(\mathbb{A}, \mathbb{B})$ there is a one-to-one mapping to the conditionals $D(\mathbb{A}|\mathbb{B})$, and the EM-like algorithm given in Sect. 7 converges to $D(\mathbb{B})$, s.t. $D(\mathbb{A}, \mathbb{B}) = D(\mathbb{A}|\mathbb{B}) \odot (I_{\mathbb{A}} \otimes D(\mathbb{B}))$.

Finally, we will reason in a simple case that generalized probability space is more “connected” and a clever algorithm might be able to exploit this. Assume zero is encoded as the distribution $(1, 0)$ and one as the distribution $(0, 1)$. Moving from the zero distribution to the one distributions can be done by lowering the probability of the first component and increasing the probability of the second. As density matrices, zero and one would be $\begin{pmatrix} 1 & 0 \\ 0 & 0 \end{pmatrix}$ and $\begin{pmatrix} 0 & 0 \\ 0 & 1 \end{pmatrix}$, respectively. Note that the eigensystem for both matrices is the identity matrix and there is now a second way to go from zero to one that keeps the eigenvalues/probabilities fixed but swaps the eigenvectors:

$$\begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix} \begin{pmatrix} 1 & 0 \\ 0 & 0 \end{pmatrix} \begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix} = \begin{pmatrix} 0 & 0 \\ 0 & 1 \end{pmatrix}.$$

Acknowledgements We are grateful to Allen Van Gelder for making us aware of Alexa (2002). Also many thanks to Torsten Ehrhardt who first proved to us the range intersection property OP1 and the \log^+ formula OP2 for the \odot operation.

Open Access This article is distributed under the terms of the Creative Commons Attribution Noncommercial License which permits any noncommercial use, distribution, and reproduction in any medium, provided the original author(s) and source are credited.

Appendix: Quantum-mechanical interpretation of conditional probability $D(a|b)$

We will now show how to interpret the conditional probability $D(a|b)$ in terms of two quantum measurements. The two measurements will be performed one after another on the joint density $D(\mathbb{A}, \mathbb{B})$ and $D(a|b)$ will be a probability of outcome 1 for the second measurement given the first measurement had outcome 1. First, we measure $D(\mathbb{A}, \mathbb{B})$ with event $I_{\mathbb{A}} \otimes bb^{\top}$. Assume that we get outcome 1. Using the generalization of collapse rule for events (see e.g. Nielsen and Chuang 2000), the successor density matrix can be computed as follows:

$$\hat{D}(\mathbb{A}, \mathbb{B}) = \frac{(I_{\mathbb{A}} \otimes bb^{\top})D(\mathbb{A}, \mathbb{B})(I_{\mathbb{A}} \otimes bb^{\top})}{\text{tr}((I_{\mathbb{A}} \otimes bb^{\top})D(\mathbb{A}, \mathbb{B})(I_{\mathbb{A}} \otimes bb^{\top}))}.$$

The second measurement consists of measuring the updated joint with event $aa^{\top} \otimes I_{\mathbb{B}}$. Now the probability for getting outcome 1 is computed as:

$$\begin{aligned} \text{tr}(\hat{D}(\mathbb{A}, \mathbb{B})(aa^{\top} \otimes I_{\mathbb{B}})) &= \frac{\text{tr}((I_{\mathbb{A}} \otimes bb^{\top})D(\mathbb{A}, \mathbb{B})(I_{\mathbb{A}} \otimes bb^{\top})(aa^{\top} \otimes I_{\mathbb{B}}))}{\text{tr}((I_{\mathbb{A}} \otimes bb^{\top})D(\mathbb{A}, \mathbb{B})(I_{\mathbb{A}} \otimes bb^{\top}))} \\ &\stackrel{\text{KP2+cycle}}{=} \frac{\text{tr}(D(\mathbb{A}, \mathbb{B})(aa^{\top} \otimes bb^{\top}))}{\text{tr}(D(\mathbb{A}, \mathbb{B})(I_{\mathbb{A}} \otimes bb^{\top}))} = \frac{D(a, b)}{\text{tr}(D(\mathbb{A}, \mathbb{B})(I_{\mathbb{A}} \otimes bb^{\top}))}. \end{aligned}$$

The denominator can be simplified using partial trace properties:

$$\text{tr}(D(\mathbb{A}, \mathbb{B})(I_{\mathbb{A}} \otimes bb^{\top})) \stackrel{\text{PT2}}{=} \text{tr}(\text{tr}_{\mathbb{A}}(D(\mathbb{A}, \mathbb{B})(I_{\mathbb{A}} \otimes bb^{\top}))) \stackrel{\text{PT3}}{=} \text{tr}(\underbrace{\text{tr}_{\mathbb{A}}(D(\mathbb{A}, \mathbb{B}))}_{D(\mathbb{B})}bb^{\top}) = D(b).$$

Therefore the probability of outcome 1 on the second measurement (given the first outcome was 1) is:

$$\text{tr}(\hat{D}(\mathbb{A}, \mathbb{B})(aa^{\top} \otimes I_{\mathbb{B}})) = \frac{D(a, b)}{D(b)} \stackrel{\text{CP4}}{=} D(a|b).$$

References

- Alexa, M. (2002). Linear combination of transformations. In *SIGGRAPH'02: Proceedings of the 29th annual conference on computer graphics and interactive techniques* (pp. 380–387). New York: ACM Press.
- Bernstein, D. S. (2005). *Matrix mathematics: theory, facts, and formulas with application to linear systems theory*. Princeton: Princeton University Press.
- Bhatia, R. (1997). *Matrix analysis*. Berlin: Springer.
- Bužek, V., Drobný, G., Derka, R., Adam, G., & Wiedemann, H. (1999). Quantum state reconstruction from incomplete data. *Chaos Solitons Fractals*, 10, 981–1074.
- Caves, C. M., Fuchs, C. A., Manne, K. K., & Renes, J. M. (2004). Gleason-type derivations of the quantum probability rule for generalized measurements. *Foundations of Physics*, 34, 193–209.
- Cerf, N. J., & Adami, C. (1999). Quantum extension of conditional probability. *Physical Review A*, 60(2), 893–897.
- Feynman, R. P. (1972). *Statistical mechanics: a set of lectures*. Reading: Addison-Wesley.
- Gleason, A. (1957). Measures on the closed subspaces of a Hilbert space. *Indiana University Mathematics Journal*, 6, 885–893.
- Holevo, A. S. (2001). *Lecture notes in physics. Monographs: Vol. 67. Statistical structure of quantum theory*. Berlin, New York: Springer.
- Kato, T. (1978). Trotter's product formula for an arbitrary pair of self-adjoint contraction semigroups. *Topics in Functional Analysis (Advances in Mathematics—Supplementary Studies)*, 3, 185–195.

- Kivinen, J., & Warmuth, M. K. (1997). Additive versus exponentiated gradient updates for linear prediction. *Information and Computation*, 132(1), 1–64.
- Kivinen, J., & Warmuth, M. K. (1999). Averaging expert predictions. In *Lecture notes in artificial intelligence: Vol. 1572. Computational learning theory, 4th European conference (EuroCOLT'99), Nordkirchen, Germany, March 29–31, 1999, Proceedings* (pp. 153–167). Berlin: Springer.
- Nielsen, M. A., & Chuang, I. L. (2000). *Quantum computation and quantum information*. Cambridge: Cambridge University Press.
- Olivares, S., & Paris, M. G. A. Quantum estimation via the minimum Kullback entropy principle. *Physical Review A*, 76, 2007.
- Schack, R., Brun, T. A., & Caves, C. M. (2001). Quantum Bayes rule. *Physical Review A*, 64, 014305.
- Simon, B. (1979). *Functional integration and quantum physics*. San Diego: Academic Press.
- Singh, R., Warmuth, M. K., Raj, B., & Lamere, P. (2003). Classification with free energy at raised temperatures. In *Proc. of EUROSPEECH 2003, September 2003* (pp. 1773–1776).
- Tsuda, K., Raätsch, G., & Warmuth, M. K. (2005). Matrix exponentiated gradient updates for on-line learning and Bregman projections. *Journal of Machine Learning Research*, 6, 995–1018.
- Warmuth, M. K. (2005). Bayes rule for density matrices. In *Advances in neural information processing systems 18 (NIPS'05)*. Cambridge: MIT Press.
- Warmuth, M. K. (2007). Winnowing subspaces. In *Proceedings of the 24th international conference on machine learning (ICML'07)*. New York: ACM.
- Warmuth, M. K., & Kuzmin, D. (2006). Online variance minimization. In *Proceedings of the 19th annual conference on learning theory (COLT'06), Pittsburg, June 2006*. New York: Springer.
- Warmuth, M. K., & Kuzmin, D. (2008). Randomized PCA algorithms with regret bounds that are logarithmic in the dimension. *Journal of Machine Learning Research*, 9, 2217–2250.
- Zellner, A. (1998). Optimal information processing and Bayes's theorem. *The American Statistician*, 42(4), 278–284.