

LINEARLY INDEPENDENT FEATURE EXTRACTION

Anonymous authors

Paper under double-blind review

ABSTRACT

We argue that domain invariance is fundamentally limiting as an objective for out-of-domain generalization (OOD) and propose a more nuanced alternative: Modeling a full spectrum of dependence on the state of the world. We make this objective tractable by developing a spectral theory, grounded in a novel operator algebra, that is formally equivalent to information-theoretic measures of dependence. The culmination is Linearly Independent Feature Extraction (LIFE): An algorithm for learning representations with controllable state-dependence, implemented using a simple eigensolver. Analytical evaluation on known data-generating processes demonstrates that LIFE recovers oracle-level features. Empirically, on linear hypothesis LIFE outperforms current gold standards and, on some datasets, even surpasses deep invariant models. A broadly applicable dynamic theory of state-dependence emerges.

1 INTRODUCTION

This is a work in representation learning (Bengio et al., 2013). The subject study is the adaptation of learned representations to large shifts in the state of the world. The objects of analysis are *stateful input streams*; Sequences of inputs whose distribution depends on a discrete random variable, hereafter *the state*.

The payoff is twofold: (i) a reframing of OOD generalization as the control and decomposition of state-dependence, and (ii) a novel spectral theory reducing this complex adaptation objective to a standard symmetric eigenproblem.

The life cycle of a learning machine follows a binary rhythm of learning and inference. During the learning phase, models are trained on *learning streams*, which are often scarce and costly to acquire. This scarcity amplifies the core challenge of out-of-domain generalization: Performing reliably on new *inference streams* while allowing for large distributional shifts.

Domain invariant learning espouses a dichotomic classification of features into domain invariant and spurious. However, persistent failures (Gulrajani & Lopez-Paz, 2020) require a direct re-evaluation of invariance itself as a primary objective for out-of-domain adaptation.

We argue for a more granular control of representations: A continuous spectrum of state-dependence that asks, how much does a feature depend on the state of the world? For this spectrum to be practical, it must be interpretable, computationally tractable, and provide sufficient conditions for alignment with the learning phase objectives during inference.

Our contribution is to construct a Gaussian theory of dependence on the state rooted in information theory and accounting for regularity on held out states, readily implemented using an eigensolver, thus offering both the necessary language and algorithmic underpinning to rigorously identify and address the fundamental challenge of out-of-domain generalization. We proceed in three steps.

First, a spectral theory of state dependence is constructed by developing a novel operator algebra, \mathcal{Q} , which unifies the treatment of both conditional and unconditional dependence and allows for

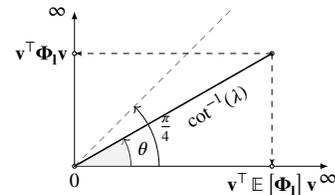


Figure 1: Spectral trigonometry of state-dependence. The dependence of an input stream, Φ_1 , and the state \mathbf{l} increases as the eigenvalue the ratio of the embedding of stream to that of its expectation with respect to the state deviates from 1.

the decomposition of representations into subspaces ranked by the strength of their coupling on the state via a standard trace-minimization argument (Fan, 1949). This will yield a simple, well defined, optimization objective for learning robust representations solvable using a single call to an eigensolver.

Second, regularity of inference is analyzed by modeling inference streams as **global perturbations** of the learning ones. The problem of ensuring *faithful processing*, the trustworthy application of learned representations to held-out inference streams, reduces to the stability of a structured symmetric eigenproblem (Demmel & Kågström, 1987) in \mathcal{Q} . Sufficient conditions are expressed in terms of the eigengap of learned representations.

Third, out-of-domain generalization is reframed. We shift the primary objective from the idealized pursuit of perfect invariance to the more practical goal of learning representations whose dependence on the state is stable across domain shifts. This provides a crucial property for deployed learning systems: Predictable out-of-domain performance that aligns with the intent established during the learning phase.

We analytically and empirically evaluate LIFE as a theory of out-of-domain generalization. Analytically, on established OOD data generating processes (Rosenfeld et al., 2020), LIFE attains oracle-level performance in both supervised and unsupervised tasks. Empirically, on linear hypothesis LIFE outperforms Invariant Risk Minimization (Arjovsky et al., 2019) while offering a simpler and stabler optimization problem, exceeds group Distributionally Robust Optimization on worst-group generalization (Sagawa et al., 2019), and, on some datasets, even outperforms deep invariant learning methods despite its linearity.

2 SPECTRAL DECOMPOSITION FOR STATE-DEPENDENCE

This section builds a spectral theory for measuring and decomposing dependence on the state. Section 2.1 introduce the fundamental data-generation and acquisition processes under variable states. Section 2.2 constructs an operator division algebra allowing for algebraic probabilistic and information theoretic calculus, leading to spectral expression of symmetrized information theoretic measure in section 2.3 and subsequently the formulation of a simple optimization objective and an algorithm for learning representations with a complete spectrum of dependence on the state in section 2.4.

First, we establish our notation as follows. We use $=$ for equality and $:=$ for assignment. Let $(\mathcal{X}, \langle \cdot, \cdot \rangle)$ be a finite dimensional inner product space. We adopt Householder’s notation (Goodfellow et al., 2017); x denotes a scalar, $\mathbf{x} = [x_1, \dots, x_m]^T$ a vector, and $\mathbf{X} = [\mathbf{x}_1 \dots \mathbf{x}_n]$ a matrix.

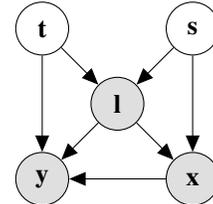
2.1 STATEFUL INPUT STREAMS

A *state* is a random vector on a discrete *state space*. A *stateful input stream* is a random variable \mathbf{x} on \mathcal{X} following a Gaussian distribution indexed by a state \mathbf{s} on a state space \mathcal{S} , $\mathbf{x} \sim \mathcal{N}(\boldsymbol{\mu}_{\mathbf{s}}, \boldsymbol{\Sigma}_{\mathbf{s}})$. A *stateful processing rule* is a conditionally Gaussian linear model (Bishop, 2011) from \mathcal{X} to \mathcal{Y} , indexed by a state \mathbf{t} on a state space \mathcal{T} , $\mathbf{x} \mapsto \mathbf{r}(\mathbf{x}) := \mathbf{A}_{\mathbf{t}}\mathbf{x} + \mathbf{b}_{\mathbf{t}} + \sqrt{\boldsymbol{\Gamma}_{\mathbf{t}}}\boldsymbol{\varepsilon}$ where $\boldsymbol{\varepsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$.

The composition of stateful input streams and processing rules defines *jointly stateful joint input streams* as defined below in **DGP-0**.

DGP-0 (Gaussian Interstate). Let \mathbf{x} be stateful input stream, \mathbf{r} a stateful processing rule, define $\mathbf{y} := \mathbf{r}(\mathbf{x})$. Then, $(\mathbf{x}, \mathbf{y}) \sim \mathcal{N}(\boldsymbol{\mu}_{[\mathbf{s}, \mathbf{t}]}, \boldsymbol{\Sigma}_{[\mathbf{s}, \mathbf{t}]})$ with,

$$\boldsymbol{\mu}_{[\mathbf{s}, \mathbf{t}]} := \begin{bmatrix} \boldsymbol{\mu}_{\mathbf{s}} \\ \mathbf{A}_{\mathbf{t}}\boldsymbol{\mu}_{\mathbf{s}} + \mathbf{b}_{\mathbf{t}} \end{bmatrix} \text{ and, } \boldsymbol{\Sigma}_{[\mathbf{s}, \mathbf{t}]} := \begin{bmatrix} \boldsymbol{\Sigma}_{\mathbf{s}} & \boldsymbol{\Sigma}_{\mathbf{s}}\mathbf{A}_{\mathbf{t}}^T \\ \mathbf{A}_{\mathbf{t}}\boldsymbol{\Sigma}_{\mathbf{s}} & \mathbf{A}_{\mathbf{t}}\boldsymbol{\Sigma}_{\mathbf{s}}\mathbf{A}_{\mathbf{t}}^T + \boldsymbol{\Gamma}_{\mathbf{t}} \end{bmatrix}.$$



While shifts in joint input streams are function of the joint state $[\mathbf{s}, \mathbf{t}]$ during data-generation. The *data acquisition process* does not observe the individual substates while being able to separate joint states. This means that there exists a state \mathbf{l} , a state space \mathcal{L} , and some unobserved bijection g such that $\mathbf{l} = g(\mathbf{s}, \mathbf{t})$. Hence, the data acquisition process stores $(\mathbf{l}, \mathbf{x}, \mathbf{y})$, see figure 2.

Figure 2: Probabilistic graphical for **DGP-0** and its data acquisition process.

Our goal is to construct an interpretable and practical spectral theory leading to measuring input streams dependence on the state and decomposing them into a discrete set of components ranked by

108 their dependence on the state. This dependence can be unconditional by considering an input stream
 109 or joint input stream as a unit, or conditional, by partitioning a joint input stream in two, in which
 110 case the measured dependence is that of a partition on the state given the remaining partition.

111 In order for the theory to be interpretable we desire it to mirror its information theoretic counterpart.
 112 Information theory of dependence is built on convex functionals of likelihood ratios and their factor-
 113 izations. For the mirroring to work, we must give an operative spectral meaning to not only input
 114 streams but also operations on input streams, in particular, their ratio and expectation with respect to
 115 the state.

116 In order for theory to be practical it should readily lead to learning representations with controllable
 117 dependence on the state through a direct and clearly stated optimization problem.

118 The foundation of our approach is to embed **DGP-0** into a division algebra of positive operators, \mathcal{Q} .
 119 This will lead to sizable conceptual and computational simplifications.

120 **Conceptually**, the algebra will allow for a unified algebraic probabilistic and information theoretic
 121 calculus; For instance, unconditional, conditional, and jointly Gaussian distributions, as well as
 122 their ratios and expectations, are all represented as single element of the operator algebra. The
 123 unification leads to disaggregation of the dependence on the state in terms of unconditional/conditional
 124 components, and subsequently, to unsupervised/supervised feature extractors.

125 **Computationally**, simultaneous diagonalization of collections of non-commuting matrices is no-
 126 toriously difficult. The algebra folds collections of non-commuting matrices into one symmetric
 127 eigenproblem, thus sidestepping the lack of Schur decomposition for more than pairs of matrices
 128 (Kressner, 2005), the simultaneous diagonalization by congruence problem (Bustamante et al., 2020),
 129 or the polynomial eigenvalue problem (Gohberg et al., 2009). Measuring dependence on the state and
 130 learning state controlled representations are reduced to standard symmetric eigenvalue problem.

133 2.2 SMOOTH OPERATOR DIVISION ALGEBRA

134 The fundamental building block of our approach is the Φ -map. This map generalizes an idea
 135 initially introduced by Siegel (1943) in the context of symplectic geometry and more recently, used
 136 in information geometry (Calvo & Oller, 1990) and manifold optimization (Hosseini & Sra, 2015).
 137 Before stating the formal definition, let's establish some intuition. A Gaussian distribution is fully
 138 described by its mean and covariance. In its simplest form, the Φ -map does embeds them jointly in a
 139 structured subset of the cone of positive definite symmetric matrix. This block matrix structure is
 140 carefully designed to turn probabilistic operations into algebraic ones. Examples of Φ -embeddings
 141 for two state bivariate Gaussian distributions are shown in figure 3.

142 More formally, writing $\mathcal{M}(m, n) = (\mathcal{M}_{m \times n}(\mathbb{R}), \langle \cdot, \cdot \rangle_F)$ the inner
 143 product space of rectangular matrices under the Frobenius norm
 144 induced by the Euclidean inner product on \mathcal{X} , $\mathcal{S}(m)$ the subspace
 145 of symmetric matrices, and $\mathcal{S}_>(m)$ the cone of positive definite
 146 matrices in $\mathcal{S}(m)$.

147 **Definition 1** (Φ -map). Given $m \geq n$ positive integers, the Φ -
 148 map is the operator valued function,

$$149 \begin{cases} \mathcal{S}_>(m) \times \mathcal{M}(m, n) & \longrightarrow \mathcal{Q}(m, n) \subset \mathcal{S}_>(m+n), \\ (\mathbf{H}, \mathbf{G}) & \longmapsto \Phi(\mathbf{H}, \mathbf{G}) := \begin{bmatrix} \mathbf{H} + \mathbf{G}\mathbf{G}^\top & \mathbf{G} \\ \mathbf{G}^\top & \mathbf{I}_n \end{bmatrix}. \end{cases}$$

150 Recall the direct sum (Singh et al., 2005) of two matrices \mathbf{X} and
 151 \mathbf{Y} forms the block diagonal matrix, $\mathbf{X} \oplus \mathbf{Y} = \begin{bmatrix} \mathbf{X} & \mathbf{0} \\ \mathbf{0} & \mathbf{Y} \end{bmatrix}$. Define
 152 the function $\mathbf{X} \mapsto \mathbf{U}(\mathbf{X}) := \begin{bmatrix} \mathbf{I}_m & \mathbf{X} \\ \mathbf{0} & \mathbf{I}_n \end{bmatrix}$. Using block Gauss-Jordan
 153 reduction (Dym, 2023), any Φ -embedding $\Phi_{\mathbf{I}'} := \Phi(\mathbf{H}_{\mathbf{I}'}, \mathbf{G}_{\mathbf{I}'})$

154 can be decomposed as the conjugation of a block diagonal matrix by an upper triangular one; $\Phi_{\mathbf{I}'} =$
 155 $\mathbf{U}(\mathbf{G}_{\mathbf{I}'}) (\mathbf{H}_{\mathbf{I}'} \oplus \mathbf{I}_n) \mathbf{U}(\mathbf{G}_{\mathbf{I}'})^\top$. This decomposition allows us to define a ratio operation on \mathcal{Q} turning it
 156 into a non-commutative operator valued analog of division on the positive reals. Moreover, \mathcal{Q} is convex
 157 and hence algebraically closed under expectation with respect to the state. This will yield considerable
 158 simplification in the upcoming analysis. Intuitively, our ratio operation can be thought of as a centering

$$142 \Phi \left(\begin{array}{c} \text{Gaussian} \end{array} \right) = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}$$

(a) State 1

$$143 \Phi \left(\begin{array}{c} \text{Gaussian} \end{array} \right) = \begin{bmatrix} 1 + \mu^2 & \rho & \mu \\ \rho & 1 & 0 \\ \mu & 0 & 1 \end{bmatrix}$$

(b) State 2

Figure 3: Examples of two states Φ -embedding. (a) embedding of standard Gaussian. (b) embedding of $\mathcal{N}\left(\begin{bmatrix} \mu \\ 0 \end{bmatrix}, \begin{bmatrix} 1 & \rho \\ \rho & 1 \end{bmatrix}\right)$, rendered with $\mu = 1.2$ and $\rho = 0.6$.

162 followed by whitening. If the two distributions are identical, this ratio is simply the identity matrix.
 163 Defining the variance of a random matrix (Muirhead, 1982) \mathbf{X} as $\mathbb{V}[\mathbf{X}] = \mathbb{E}[\mathbf{X}\mathbf{X}^\top] - \mathbb{E}[\mathbf{X}]\mathbb{E}[\mathbf{X}]^\top$, the
 164 discussion above is encapsulated in theorem 1.

165 **Theorem 1** (Operator division algebra). *Let $(\mathbf{H}_1, \mathbf{G}_1) \in \mathcal{S}_>(m) \times \mathcal{M}_{m \times n}(\mathbb{R})$, and $\Phi_1 := \Phi(\mathbf{H}_1, \mathbf{G}_1)$.
 166 Then $\mathcal{Q}(m, n)$ is convex and closed under the binary ratio operations:*

167 $(\Phi_1, \Phi_{1'}) \mapsto \Phi_1 : \Phi_{1'} := \left(\mathbf{H}_{1'}^{-\frac{1}{2}} \oplus \mathbf{I}_n \right) \mathbf{U}(-\mathbf{G}_{1'}) \Phi_1 \mathbf{U}(-\mathbf{G}_{1'})^\top \left(\mathbf{H}_{1'}^{-\frac{1}{2}} \oplus \mathbf{I}_n \right)$. Moreover,

169 (i) $\Phi_1 : \mathbf{I} = \Phi_1$,

171 (ii) $\Phi_1 : \Phi_{1'} = \Phi(\mathbf{Q}_{11'}, \Delta_{11'})$, with $\mathbf{Q}_{11'} := \mathbf{H}_{1'}^{-\frac{1}{2}} \mathbf{H}_1 \mathbf{H}_{1'}^{-\frac{1}{2}}$ and $\Delta_{11'} := \mathbf{H}_{1'}^{-\frac{1}{2}} (\mathbf{G}_1 - \mathbf{G}_{1'})$,

172 (iii) $\mathbb{E}[\Phi_1 : \Phi_{1'}] = \Phi(\mathbb{E}[\mathbf{Q}_{11'}] + \mathbb{V}[\Delta_{11'}], \mathbb{E}[\Delta_{11'}])$.

174 In the following $\lambda(\mathbf{X})$ denotes the set of eigenvalues of \mathbf{X} , $\lambda^\uparrow(\mathbf{X})$ the vector of eigenvalues, including
 175 multiplicity, arrayed in non-decreasing order.

176 \mathcal{Q} is a subset of the space of symmetric matrices; The spectral theorem applies (Hilbert, 1989). Hence,
 177 by theorem 1 \mathcal{Q} carries a structured eigenvalue problem inherited by all of its elements. In fact, each
 178 element of \mathcal{Q} associates a symmetric eigenvalue problem to a rational eigenvalue problem (Xi & Saad,
 179 2015; Betcke et al., 2013); The spectrum of the former is equal to the union of the spectrum and
 180 the poles of the latter. Thankfully, it is generally simpler to solve a symmetric eigenproblem on \mathcal{Q}
 181 than a rational eigenproblem (Su & Bai, 2011). By expressing the eigenvectors/spectrum of a typical
 182 element of \mathcal{Q} as the critical points/values of the Rayleigh quotient (Sun, 1991), block coordinates
 183 optimization (Lange, 2016) explicitly characterizes this association.

184 **Lemma 1.** *Let $\mathbf{H} \in \mathcal{S}_>(m)$, $\mathbf{G} \in \mathcal{M}_{m \times n}(\mathbb{R})$, $\Phi := \Phi(\mathbf{H}, \mathbf{G}) \in \mathcal{Q}(m, n)$, and $\mathbf{w}^\top = [\mathbf{w}_1^\top \quad \mathbf{w}_2^\top]$ be a
 185 unit normalized eigenvector of Φ . Then any eigenpair (λ, \mathbf{W}) of Φ verifies,*

187 (i) *If $\lambda = 1$, $\mathbf{G}^\top \mathbf{w}_1 = 0$ and $(\mathbf{H} - \mathbf{I})\mathbf{w}_1 = -\mathbf{G}\mathbf{w}_2$.*

188 (ii) *If $\lambda \neq 1$, (λ, \mathbf{w}) is an eigenpair of Φ if and only (λ, \mathbf{w}_1) is an eigenpair of the rational
 189 eigenvalue problem, $\mathbf{H}\mathbf{w}_1 = \lambda \left(\mathbf{I} + \frac{1}{1-\lambda} \mathbf{G}\mathbf{G}^\top \right) \mathbf{w}_1$, and $\mathbf{G}^\top \mathbf{w}_1 = (1 - \lambda)\mathbf{w}_2$.*

191 An immediate application is that lemma 1 provides a unified operative interpretation for the canonical
 192 embeddings of definition 2, as well as their ratios, and expected ratios by considering them as an
 193 element of \mathcal{Q} . The spectrum of Φ -embeddings of a Gaussian distribution can be viewed as an absolute,
 194 or magnitude, spectrum; that of a ratio as a relative, or phase, spectrum. In particular, for a the
 195 embedding of Gaussian, $\Phi := \Phi(\Sigma, \mu)$ any eigenpair $(\lambda, \mathbf{w}^\top := [\mathbf{w}_1^\top \quad w_2])$ with $\|\mathbf{w}\|_2 = 1$ verifies
 196 $\lambda(\Phi) = \mathbf{w}_1^\top \Sigma \mathbf{w}_1 + (\langle \mathbf{w}_1, \mu \rangle + w_2)^2$. When $\mu = \mathbf{0}$ any eigenpair of Σ is one of Φ . When $\mu \neq \mathbf{0}$ the
 197 spectrum of Φ differ from that of Σ through the addition of a quadratic term accounting for the mean.
 198 In contrast to the spectrum of Σ , that of Φ is mean aware. For a ratio of unconditional canonical embed-
 199 dings, any such eigenpair verifies $\lambda(\Phi_1 : \Phi_2) = \mathbf{w}_1^\top (\Sigma_2^{-\frac{1}{2}} \Sigma_1 \Sigma_2^{-\frac{1}{2}}) \mathbf{w}_1 + \left(\left\langle \mathbf{w}_1, \Sigma_2^{\frac{1}{2}} (\mu_1 - \mu_2) \right\rangle + w_2 \right)^2$.

201 Each eigenvalue separates into two different components: One accounting for variations of the means
 202 and another for that of the covariances. The spectrum is said to be relative as it quantifies how Φ_1
 203 and Φ_2 differ. This “phase” information is exactly what we need to measure dependence on the state.

204 205 206 207 2.3 A SPECTRAL THEORY OF STATE-DEPENDENCE

208 In this section we spectrally express symmetric information-theoretic measure of dependence on the
 209 state by writing them as the trace of single element of \mathcal{Q} ; State dependence is read from the spectrum
 210 of a single symmetric operator, preparing the optimization in section 2.4.

211 The mutual and conditional mutual information (Shannon, 1948) exhibit properties that render them
 212 well-suited for assessing dependence: They vanish if and only if the variables are independent and
 213 adhere to a chain rule (Polyanskiy & Wu, 2024) that disaggregates joint dependence into marginal
 214 and conditional components. When the joint distribution is absolutely continuous with respect
 215 to the product of marginals (Belghazi et al., 2018), the mutual information corresponds to the
 Kullback–Leibler divergence (Kullback, 1968).

In this study, we employ the symmetrized version of the KL divergence to define symmetrized mutual information. Recall that if \mathbb{P} is absolutely continuous with respect to \mathbb{Q} , the Jeffreys divergence (Csiszár, 1972) is defined as $D_J(\mathbb{P} \parallel \mathbb{Q}) := KL(\mathbb{P} \parallel \mathbb{Q}) + KL(\mathbb{Q} \parallel \mathbb{P})$. Unlike the asymmetric KL divergence, Jeffreys' divergence does not include a log-partition function; This will enable exact spectral expressions over \mathcal{Q} .

In the probabilistic graphical model of figure 2, we define the *unconditional dependence on the state* as $S(\mathbf{x}; \mathbf{I}) := D_J(\mathbb{P}_{\mathbf{xI}} \parallel \mathbb{P}_{\mathbf{x}} \otimes \mathbb{P}_{\mathbf{I}})$, and the *conditional dependence on the state* $S(\mathbf{y}; \mathbf{I} \mid \mathbf{x}) := D_J(\mathbb{P}_{\mathbf{yI} \mid \mathbf{x}} \parallel \mathbb{P}_{\mathbf{y} \mid \mathbf{x}} \otimes \mathbb{P}_{\mathbf{I} \mid \mathbf{x}} \mid \mathbb{P}_{\mathbf{x}})$. First, we must verify that the symmetrized mutual information are as equally suited to measuring dependence as their unsymmetric counterparts.

Proposition 1 (Behaviour after measurement). *If $(\mathbf{I}, \mathbf{x}, \mathbf{y})$ are distributed according to figure 2. Then the following properties hold,*

CALIBRATED INDEPENDENCE: $S(\mathbf{x} \mid \mathbf{I}) = 0 \iff \mathbf{x} \perp\!\!\!\perp \mathbf{I}$ and $S(\mathbf{y}; \mathbf{I} \mid \mathbf{x}) = 0 \iff \mathbf{y} \mid \mathbf{x} \perp\!\!\!\perp \mathbf{I} \mid \mathbf{x}$.

CHAIN RULE OF INFORMATION: $S(\mathbf{x}; \mathbf{y} \mid \mathbf{I}) = S(\mathbf{x}; \mathbf{I}) + S(\mathbf{y}; \mathbf{I} \mid \mathbf{x})$.

STATE INDEPENDENT BAYES: $S(\mathbf{x}; \mathbf{I}) = S(\mathbf{y}; \mathbf{I} \mid \mathbf{x}) = 0 \implies S(\mathbf{y}; \mathbf{I}) = S(\mathbf{x}; \mathbf{I} \mid \mathbf{y}) = 0$.

We now leverage the algebraic structure offered by \mathcal{Q} to uniformly handle unconditional and conditional dependence on the state. As information measure of dependence can be expressed as functionals of likelihood ratios. It is enough for us to represent marginal, conditional and joint distribution as element of \mathcal{Q} ; Their ratios, and expectations, will be automatically handled by the ratio operation on \mathcal{Q} . This leads to the definition of the following canonical elements.

Definition 2 (Canonical Φ -embeddings). In the notation of DGP-0, a Φ -embedding $\Phi(\mathbf{H}_1, \mathbf{G}_1)$ is called *canonical of the unconditional kind* if, $\mathbf{H}_1 = \Sigma_1$ and $\mathbf{G}_1 = \mu_1$, or of the *conditional kind* if $\mathbf{H}_1 = \Gamma_1$ and $\mathbf{G}_1 = [\mathbf{A}_1 \quad \mathbf{b}_1] \in [\Phi(\Sigma_1, \mu_1)]^{\frac{1}{2}}$.

Note that, the conditional canonical Φ -embedding demonstrates that \mathcal{Q} can represent more than distributions and their ratios.

Now that we have our algebraic tools, we can define our measure of state-dependence. The core idea is to measure how much each state-specific distribution Φ_1 deviates from a single, fixed reference element of \mathcal{Q} . The most natural reference is the expectation with respect to the state, the center of mass $\mathbb{E}[\Phi_1]$. The relative spectrum of the random ratio $\mathbb{E}[\Phi_1] : \Phi_1$ quantities how much Φ_1 departs from the state independent $\mathbb{E}[\Phi_1]$, see figure 1. Just as the aggregation of the log likelihood ratio of the joint to the product of the marginal characterizes the mutual information, $\mathbb{E}[\mathbb{E}[\Phi_1] : \Phi_1]$ characterize the dependence on the state. More formally,

Theorem 2 (Spectral dependence). $\frac{1}{2} \text{Tr}[\mathbb{E}[\mathbb{E}[\Phi_1] : \Phi_1] - \mathbf{I}]$ is equal to $S(\mathbf{x}; \mathbf{I})$, if Φ is a canonical Φ -embedding of the unconditional kind, or $S(\mathbf{y}; \mathbf{I} \mid \mathbf{x})$, if Φ is of the conditional kind.

2.4 LINEARLY INDEPENDENT FEATURE EXTRACTION

We aim to identify subspaces over which state dependences decompose. Our objective is to determine the optimal rank- k subspace that minimizes these dependences. Given the symmetry of $\mathbb{E}[\mathbb{E}[\Phi_1] : \Phi_1]$, the Ky-Fan trace minimization principle (Fan, 1949) readily applies. Writing $\text{Eig}(\mathbf{X}, \lambda) = \text{span}\{\mathbf{w} : \mathbf{X}\mathbf{w} = \lambda\mathbf{w}\}$ as the eigenspace of \mathbf{X} associated to λ , we have,

Corollary 1. *Let Φ_1 be a state index canonical Φ -embedding. Then,*

$$\mathbb{R}^{m+n} = \bigoplus_{\lambda \in \lambda(\mathbb{E}[\mathbb{E}[\Phi_1] : \Phi_1])} \text{Eig}(\mathbb{E}[\mathbb{E}[\Phi_1] : \Phi_1], \lambda)$$

$$\min_{\mathbf{W}^T \mathbf{W} = \mathbf{I}_k} \text{Tr}[\mathbf{W}^T \mathbb{E}[\mathbb{E}[\Phi_1] : \Phi_1] \mathbf{W}] = \sum_{i=1}^k \lambda_i^{\uparrow}(\mathbb{E}[\mathbb{E}[\Phi_1] : \Phi_1]).$$

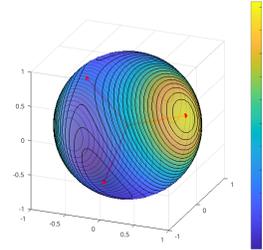


Figure 4: Contour of the expected ratio of the two state model of figure 3. Eigenvectors and eigenvalues of the value function defined in corollary 1 are show in red.

Corollary 1 defines two fundamental algebraically invariant subspaces. 1) **The state independent**

```

270 Algorithm 1 LIFE: A procedure
271 Require:  $\mathcal{D} = [\mathcal{D}_1, \dots, \mathcal{D}_{|\mathcal{L}|}]$  # Learning data.
272 Require:  $\mathbf{p} = [p_1, \dots, p_{|\mathcal{L}|}]$  # priors.
273 Require:  $k \in \mathbb{N}$  # Number of eigenvectors.
274 Require:  $\text{mode} \in \{\text{uncond}, \text{cond}\}$ .
275 Ensure: A feature extractor  $\mathbf{W}$ .
276
277 1 procedure LIFE( $\mathcal{D}, \mathbf{p}, k, \text{mode}$ )
278 2   for  $l \leftarrow 1$  to  $|\mathcal{L}|$  do
279 3      $\Phi[l] \leftarrow \text{EstimatePhi}(\mathcal{D}[l], \text{mode})$ 
280 4   end for
281 5    $\underline{\Phi} \leftarrow \text{FormExpectedRatio}(\Phi, \mathbf{p})$ 
282 6    $(\Lambda, \mathbf{W}) \leftarrow \text{Eig}(\underline{\Phi})$ 
283 7   return  $\mathbf{W}[:, : k]$ , #  $k$ -smallest eigenvectors.
284 8 end procedure

```

eigenspace, $\mathcal{E}_\perp := \text{Eig}(\mathbb{E}[\mathbb{E}[\Phi_1] : \Phi_1], 1)$, if it exists, corresponds to the space of zero dependence on state. 2) **The rank- k subspace of minimal dependence** defined as the range of the spectral projector, $\mathcal{E}_k := \mathcal{R}(\mathbf{W}\mathbf{W}^\top)$.

3 REGULARITY OF INFERENCE

This section considers question of inference pertaining to dependence on the state. We establish sufficient conditions under which state-independent and rank- k minimal dependence subspaces can faithfully process inference streams. Representation spaces hold semantic meaning (Zeiler & Fergus, 2013); Thus, it is crucial to ensure that this meaning aligns consistently with the inference streams.

We first show how the expected ratio of inference streams can be represented as **global** perturbations of the learning one. Next, we derive sufficient conditions for faithful processing in terms of the magnitude of this perturbation. Finally, we demonstrate that, under these conditions, the dependence on the state remains regular and predictable throughout inference.

3.1 INFERENCE STREAMS AS GLOBAL PERTURBATION

We start by partitioning the state space into disjoint learning and inference states $\mathcal{L} = \underline{\mathcal{L}} \sqcup \overline{\mathcal{L}}$. Let $\underline{\Phi} := \mathbb{E}[\mathbb{E}[\Phi_1 | \underline{\mathcal{L}}] : \Phi_1 | \underline{\mathcal{L}}]$ be the expected ratio over the learning states. Similarly, let $\overline{\Phi} = \mathbb{E}[\mathbb{E}[\Phi_1 | \overline{\mathcal{L}}] : \Phi_1 | \overline{\mathcal{L}}]$, be the expected ratio of inference streams which we can also write as $\overline{\Phi} = \underline{\Phi} + \mathbf{E}$. The fundamental question is: When is an algebraically invariant space of the expected ratio of the learning streams also one for inference streams? We will reason on the norm of the perturbation as it can cause eigenvalue coalescence (Demmel, 1986); Where at least one eigenvalue shifts from a chosen eigencluster to its spectrum's complement.

3.2 FAITHFUL PROCESSING AND REGULARITY OF DEPENDENCE

According to the Schur decomposition theorem (Mayers et al., 1986), identifying conditions for the spectral resolution of the perturbation to be block upper triangular is sufficient (Kressner, 2005). This occurs if and only if an algebraic Riccati matrix equation (Lancaster & Rodman, 1995) has a solution. Conditions ensuring this solution are based on norms of perturbation blocks and spectrum eigencluster separation (Stewart et al., 1990). Specifically, if the algebraically invariant subspaces correspond to the first k smallest eigenvalues without multiplicity, the separation from the spectrum is the eigengap, $\lambda_{k+1}^\uparrow(\underline{\Phi}) - \lambda_k^\uparrow(\underline{\Phi})$, as detailed in the following proposition.

Theorem 3 (Conditions for faithful processing). *Let \mathcal{E}_k be the rank- k subspace of minimal dependence of $\underline{\Phi}$. Let $\overline{\Phi} = \underline{\Phi} + \mathbf{E}$. If $\|\mathbf{E}\|_F < \frac{1}{2} \left(\lambda_{k+1}^\uparrow(\underline{\Phi}) - \lambda_k^\uparrow(\underline{\Phi}) \right)$, then \mathcal{E}_k is an algebraically invariant subspace of $\overline{\Phi}$.*

Now that we have conditions for faithful processing we can control the regularity of the dependence on the state.

Corollary 2 (Regularity of dependence on the state). *Let \mathbf{I}' be an independent copy of \mathbf{I} , $p := \mathbb{P}((\mathbf{I}, \mathbf{I}') \in \underline{\mathcal{L}} \times \underline{\mathcal{L}})$, \mathbf{W}_k be a matrix with k orthonormal columns such that $\mathcal{R}(\mathbf{W}_k) = \mathcal{E}_k$. Then, under conditions for faithful processing,*

$$\mathrm{Tr}[\mathbf{W}_k^\top \mathbb{E}[\mathbb{E}[\Phi_1] : \Phi_1] \mathbf{W}_k] \leq \mathrm{Tr}[\mathbf{W}_k^\top \underline{\Phi} \mathbf{W}_k] + \frac{(1-p)}{2} (\lambda_{k+1}^\uparrow(\Phi) - \lambda_k^\uparrow(\Phi)).$$

4 RELATED WORK

Modal decompositions are foundational to correspondence analysis (Hirschfeld, 1935) and Lancaster’s distribution theory (Lancaster, 1958), finding use in strong data processing inequalities (Polyanskiy & Wu, 2016; Raginsky, 2014) and representation learning (Huang et al., 2024). They achieve a spectral decomposition of bivariate distribution dependencies (Lancaster, 1958) via singular value decomposition (Schmidt, 1907) (SVD) of the conditional expectation operator (Makur, 2019). However, the inherent dependence on SVD limits their usefulness as the number of component quantity is at most equal to the number of states which makes them unsuited for our setting.

Relative Matrix Decomposition The Generalized Singular Value Decomposition (GSVD) (Van Loan, 1976; Paige & Saunders, 1981) jointly decomposes two matrices with identical column counts into three bases: two orthonormal bases specific to each matrix and one shared basis. The GSVD finds applications in prenatal EKG (Callaerts et al., 1990) and genetics (Aiello et al., 2018). Unfortunately, the GSVD operates only on two matrices and not at the distributional level. Extensions allowing decomposition of more than two matrices have been proposed (De Lathauwer, 2011; Ponnappalli et al., 2011; Khamidullina et al., 2020), however these approach lose either the orthogonality of the subspaces or the exactness of the decomposition (Ponnappalli et al., 2011).

Spectral Perspectives on Representation Learning A longstanding approach interprets neural representations spectrally. Roweis & Brody (1999) examines learned linear neural representations in regression via the SVD. Linsker (1988); Deco & Obradovic (1996) relate InfoMAX to PCA via Hebbian learning. Baldi & Hornik (1989); Zhou & Liang (2017) analyze linear neural networks critical points and, using centered Gaussian distributions, and express them as permutation of principal singular vectors. Recent work extends these approach to self-supervised learning (Balestriero & LeCun, 2022) and infinite width Gaussian processes approximation of neural networks (Jacot et al., 2018; Yang, 2019).

Domain invariant learning for Out-of-Domain Generalization aims to learn domain-invariant predictors (Arjovsky et al., 2019; Mahajan et al., 2020; Krueger et al., 2020; Ye et al., 2021; Lin et al., 2022). However, empirical performance has often fallen short (Gulrajani & Lopez-Paz, 2020). Several studies (Rosenfeld et al., 2020; Kaur et al., 2022; Ahuja et al., 2020b) propose structural data-generating processes for learning invariant representations. Linear analyses of invariance are discussed in Wang & Veitch (2022); Zhang et al. (2025). In Krueger et al. (2020), invariance denotes statistical relationships conserved across domains, while in Li et al. (2020), invariant representations are defined as statistically independent of the domain. This dependency is analyzed in terms of unconditional/conditional mutual information in Tachet et al. (2020); Dong et al. (2024). Invariant Risk Minimization (IRM) (Arjovsky et al., 2019; Ahuja et al., 2020a; Lin et al., 2022; Ahuja et al., 2021; Zhou et al., 2022), which claims to estimates invariant causal predictors, remains the prevalent approach.

5 EVALUATION AS A THEORY OF OUT-OF-DOMAIN GENERALIZATION

We begin by highlighting methodological limitations in the dominant paradigm: estimating invariant causal predictors via Invariant Risk Minimization (IRM). Then we introduce LIFE as an objective for out-of-domain generalization and evaluate it both analytically and empirically.

5.1 RE-EVALUATING INVARIANCE AS A PRIMARY OBJECTIVE FOR OUT-OF-DOMAIN GENERALIZATION

The dominant out-of-distribution paradigm fixates on estimating perfectly invariant features, be it causal (Arjovsky et al., 2019) or statistical (Krueger et al., 2020), and face both methodological and practical limits. Methodologically, it treats the existence of invariant features as a precondition for learning, rather than as a falsifiable hypothesis to be tested. Practically, even when they exist and are successfully estimated, these features may lack sufficient predictive power—a failure mode known as excessive invariance (Jacobsen et al., 2018; Geirhos et al., 2020). This work argues that the fundamental challenge is one of predictable performance, where a deployed model’s behavior aligns with expectations established during learning. Therefore, **to ask only if a feature is invariant is to mistake a practical problem of degree for a philosophical problem of kind.**

5.2 ANALYTICAL EVALUATION

In this section we will demonstrate, at least in the verifiable setting established by influential OOD data-generating processes (Rosenfeld et al., 2020), that the search for causal invariant predictors reduces to identification of the state independent subspace. The appeal to the language of causality, in these analyzable settings, is an unnecessary complication. The features sought correspond precisely to a single point on our spectrum. Following, and slightly extending (Rosenfeld et al., 2020; Zhang et al., 2025), we will define DGPs at three level: 1) The topology of the conditional probabilistic graphical model, 2) The functional dependence of the nodes conditional moments on their parents, and 3) The distributions of the functional form of the conditional distribution of the nodes on their parents. The most fundamental assumption behind domain invariance is one of separability; The moments of \mathbf{x} can be written as function of 1) An *invariant* or *common* random variable \mathbf{c} that is independent of \mathbf{I} , and 2) a *spurious* or *peculiar* random variable \mathbf{s} that is dependent on \mathbf{I} . We illustrate the decomposition in figure 5. Taking $\mathbf{A} \in \mathbb{R}^{d_c \times d}$, $\mathbf{B} \in \mathbb{R}^{d_l \times d}$ and $1 \leq d_l, d_c \leq d$. Writing, $\mathbf{z}^\top = [\mathbf{c}^\top \ \mathbf{s}^\top]$ and $\mathbf{F}^\top = [\mathbf{A}^\top \ \mathbf{B}^\top]$. Define, $\mathbf{x} = \mathbf{F}^\top \mathbf{z}$. Following (Rosenfeld et al., 2020), we assume that the latent factors of variations \mathbf{z} can be recovered after mixing; Or more succinctly \mathbf{F} to be injective.

DGP-U (Unsupervised).

$$\mathbf{c} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma}) \text{ and, } \mathbf{s} \mid \mathbf{I} \sim \mathcal{N}(\boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1), \\ \mathbf{x} \mid \mathbf{I} \sim \mathcal{N}(\mathbf{A}^\top \boldsymbol{\mu} + \mathbf{B}^\top \boldsymbol{\mu}_1 \mid \mathbf{A}^\top \boldsymbol{\Sigma} \mathbf{A} + \mathbf{B}^\top \boldsymbol{\Sigma}_1 \mathbf{B}).$$

DGP-S (Supervised).

$$\mathbf{c} \mid y \sim \mathcal{N}(y\boldsymbol{\mu}, \boldsymbol{\Sigma}) \text{ and, } \mathbf{s} \mid y, \mathbf{I} \sim \mathcal{N}(y\boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1), \\ \mathbf{x} \mid y, \mathbf{I} \sim \mathcal{N}(y(\mathbf{A}^\top \boldsymbol{\mu} + \mathbf{B}^\top \boldsymbol{\mu}_1), \mathbf{A}^\top \boldsymbol{\Sigma} \mathbf{A} + \mathbf{B}^\top \boldsymbol{\Sigma}_1 \mathbf{B}).$$

For both data-generating processes the oracle distribution is defined by setting the matrix \mathbf{B} to zero, the learning states are set the states 1 and 2. For **DGP-S**, we take $y \sim \text{Rademecher}(1/2)$.

We analyze **DGP-U** and **DGP-S** through their Φ -embeddings. We form their respective canonical expected ratio; Unconditional for the former and conditional for the latter. Analysis of the spectra of the respective expected ratios reveals the existence of state independent subspaces. Moreover, the range of the spectral projector associated with the state independent subspace is precisely equal to the feature space spanned by the invariant, or oracle, component of the data-generating processes and annihilates the spurious component.

Proposition 2. *The expected ratio $\mathbb{E}[\mathbb{E}[\Phi_1] : \Phi_1]$ of both **DGP-U** and **DGP-S** admit a state independent subspace \mathcal{E}_\perp of dimension d_c . Moreover, a matrix $\mathbf{W}_1^\top = [\mathbf{W}_{11}^\top \ \mathbf{W}_{12}^\top]$ such that $\mathcal{R}(\mathbf{W}_{11}) = \mathcal{E}_\perp$. Then the orthogonal projection \mathbf{P} onto $\mathcal{R}(\mathbf{W}_{11})$ verifies,*

- (i) $\mathbf{B}\mathbf{P} = \mathbf{0}$ (Annihilation of state dependent component),
- (ii) $\mathcal{R}(\mathbf{A}^\top) = \mathcal{R}(\mathbf{P})$ (Preservation of oracle component),
- (iii) $\mathbf{P}\mathbf{x} = \mathbf{A}^\top \mathbf{c}$ for **DGP-U**, and, $\mathbf{P}\mathbf{x} = y\mathbf{A}^\top \mathbf{c}$ for **DGP-S**.

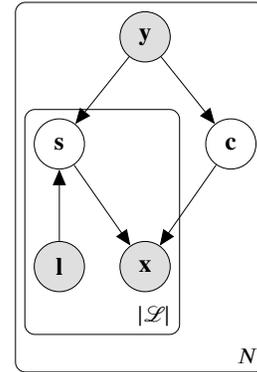


Figure 5: Probabilistic graphical model for **DGP-S**.

5.3 EMPIRICAL EVALUATION

Our empirical evaluation assesses the practical utility of LIFE across a suite of standard out-of-domain (OOD) benchmarks. We aim to determine if the extracted features (i) generalize out-of-domain, (ii) achieve robust worst-group generalization, and (iii) mitigate shortcut learning, while remaining

432 competitive against existing methods. For each experiment, we apply LIFE to extract a rank- k
 433 subspace of minimal state dependence by selecting the eigenvectors corresponding to the smallest
 434 eigenvalues of the expected ratio matrix, as per corollary 1 .

435 For all experiments on non-synthetic dataset we filter the in-
 436 puts using best k -rank approximation using randomized singu-
 437 lar value decomposition (Halko et al., 2009; Martinsson et al.,
 438 2011). Classification tasks are resolved using a logistic regres-
 439 sion (Cox, 1958) trained using L-BFGS (Liu & Nocedal, 1989).
 440

441 **Out-of-domain generalization:** ColorMNIST modifies
 442 MNIST (LeCun, 1998) by introducing color as an additional
 443 feature. Digits less than 5 are assigned to class 0; digits 5 and
 444 above to class 1. The environment variable sets the proportion of
 445 each class rendered in red or blue, creating strong dependencies
 446 between color and label that does not carry across environments.
 447 Following the exact protocol of (Arjovsky et al., 2019), results
 448 are presented in table 1.

449 **Worst-group accuracy:** Waterbirds evaluates ro-
 450 bustness to background-label associations, with
 451 images of waterbirds and land birds placed on
 452 both matching and mismatched (rare) backgrounds.
 453 Worst-group accuracy measures performance on
 454 these challenging mismatched cases, following
 455 (Sagawa et al., 2019), with results in table 2.
 456

457 **Shortcut learning:** MNIST-CIFAR (Shah et al., 2020) forms
 458 two classes by vertically concatenating MNIST digits and
 459 CIFAR-10 (Hinton, 2007) images: class 0 pairs digit zero with
 460 automobiles; class 1 pairs digit one with trucks. In this setting,
 461 MNIST digits provide a shortcut signal that is easier to classify
 462 than the associated CIFAR images. We follow the protocol of
 463 (Shah et al., 2020), with results summarized in table 3.
 464

465 6 CONCLUSION AND PERSPECTIVES

466 We introduced a spectral theory of state-dependence and its algo-
 467 rithmic realization in LIFE, demonstrating its effectiveness on
 468 the challenge of out-of-distribution generalization. The primary
 469 contribution is the computational theory itself; OOD being its
 470 first application. It offers a principled, tractable lens for analyz-
 471 ing learning systems subject to changing states of the world, e.g.
 472 continual learning, federated learning, and anomaly detection.
 473 These are not extensions by analogy but consequences of the
 474 same formulation.
 475
 476
 477
 478
 479
 480
 481
 482
 483
 484
 485

Table 1: On linear hypothesis, LIFE outperforms established baselines in generalizing to an environment with a shifted color-label correlation.

Algorithm	Accuracy
IRMv1	64.85
Spectral Decoupling	63.67
LIFE	65.68
Oracle	67.44

Table 2: By operating on ResNet-18 features, LIFE improves worst-group accuracy on the Waterbirds dataset by effectively isolating bird features from spurious background signals

Algorithm	Avg	Worst
GroupDRO	91.13	77.57
ISR-Cov	90.46 \pm 0.80	82.46 \pm 0.55
LIFE	89.38 \pm 0.33	84.17 \pm 0.67

Table 3: Even when constrained to linear features, LIFE approaches oracle-level performance on MNIST-CIFAR, successfully ignoring the shortcut feature and outperforming several deep invariant learning methods.

Algorithm	Accuracy
ERM	39.5 \pm 0.4
IRMGAME	46.7 \pm 2.1
DILU	50.2 \pm 1.7
IRMv1	51.3 \pm 3.0
REx	50.1 \pm 2.2
InvRat	52.3 \pm 0.9
BIRM (ResNet-18)	59.3 \pm 2.3
Oracle (ResNet-18)	83.5 \pm 1.5
Oracle (Linear)	63.72 \pm 0.3
LIFE (Linear)	63.5 \pm 0.35

486
487
488
489
490
491
492
493
494
495
496
497
498
499
500
501
502
503
504
505
506
507
508
509
510
511
512
513
514
515
516
517
518
519
520
521
522
523
524
525
526
527
528
529
530
531
532
533
534
535
536
537
538
539

REFERENCES

- Kartik Ahuja, Karthikeyan Shanmugam, Kush R. Varshney, and Amit Dhurandhar. Invariant Risk Minimization Games. feb 2020a. doi:[10.48550/arXiv.2002.04692](https://doi.org/10.48550/arXiv.2002.04692). (cited on page 7)
- Kartik Ahuja, Jun Wang, Amit Dhurandhar, Karthikeyan Shanmugam, and Kush R. Varshney. Empirical Or Invariant Risk Minimization? A Sample Complexity Perspective. oct 2020b. doi:[10.48550/arXiv.2010.16412](https://doi.org/10.48550/arXiv.2010.16412). (cited on page 7)
- Kartik Ahuja, Ethan Caballero, Dinghuai Zhang, Jean-Christophe Gagnon-Audet, Yoshua Bengio, Ioannis Mitliagkas, and Irina Rish. Invariance Principle Meets Information Bottleneck For Out-Of-Distribution Generalization. jun 2021. doi:[10.48550/arXiv.2106.06607](https://doi.org/10.48550/arXiv.2106.06607). (cited on page 7)
- Katherine A. Aiello, Sri Priya Ponnappalli, and Orly Alter. Mathematically Universal And Biologically Consistent Astrocytoma Genotype Encodes For Transformation And Predicts Survival Phenotype. *Apl Bioengineering*, 2(3):31909, sep 2018. doi:[10.1063/1.5037882](https://doi.org/10.1063/1.5037882). (cited on page 7)
- Martin Arjovsky, Léon Bottou, Ishaan Gulrajani, and David Lopez-Paz. Invariant Risk Minimization. jul 2019. doi:[10.48550/arXiv.1907.02893](https://doi.org/10.48550/arXiv.1907.02893). (cited on pages 2, 7, 8, and 9)
- Pierre Baldi and Kurt Hornik. Neural Networks And Principal Component Analysis: Learning From Examples Without Local Minima. *Neural Networks*, 2(1):53–58, jan 1989. doi:[10.1016/0893-6080\(89\)90014-2](https://doi.org/10.1016/0893-6080(89)90014-2). (cited on page 7)
- Randall Balestriero and Yann LeCun. Contrastive And Non-Contrastive Self-Supervised Learning Recover Global And Local Spectral Embedding Methods. may 2022. doi:[10.48550/arXiv.2205.11508](https://doi.org/10.48550/arXiv.2205.11508). (cited on page 7)
- Mohamed Ishmael Belghazi, Aristide Baratin, Sai Rajeswar, Sherjil Ozair, Yoshua Bengio, Aaron Courville, and R Devon Hjelm. Mine: Mutual Information Neural Estimation. *Arxiv*, jan 2018. doi:[10.48550/arXiv.1801.04062](https://doi.org/10.48550/arXiv.1801.04062). (cited on page 4)
- Y. Bengio, A. Courville, and P. Vincent. Representation Learning: A Review And New Perspectives. *Ieee Transactions On Pattern Analysis And Machine Intelligence*, 35(8):1798–1828, aug 2013. doi:[10.1109/tpami.2013.50](https://doi.org/10.1109/tpami.2013.50). (cited on page 1)
- Timo Betcke, Nicholas J. Higham, Volker Mehrmann, Christian Schröder, and Françoise Tisseur. Nlevp. *Acm Transactions On Mathematical Software*, 39(2):1–28, feb 2013. doi:[10.1145/2427023.2427024](https://doi.org/10.1145/2427023.2427024). (cited on page 4)
- Christopher M. Bishop. *Pattern Recognition And Machine Learning*. Springer New York, 2011. (cited on page 2)
- John Parker Burg. The Relationship Between Maximum Entropy Spectra And Maximum Likelihood Spectra. *Geophysics*, 37(2):375–376, apr 1972. doi:[10.1190/1.1440265](https://doi.org/10.1190/1.1440265). (cited on page 26)
- Miguel D. Bustamante, Pauline Mellon, and M. Victoria Velasco. Solving The Problem Of Simultaneous Diagonalization Of Complex Symmetric Matrices Via Congruence. *Siam Journal On Matrix Analysis And Applications*, 41(4):1616–1629, jan 2020. doi:[10.1137/19m1280430](https://doi.org/10.1137/19m1280430). (cited on page 3)
- D. Callaerts, B. De Moor, J. Vandewalle, W. Sansen, G. Vantrappen, and J. Janssens. Comparison Of Svd Methods To Extract The Foetal Electrocardiogram From Cutaneous Electrode Signals. *Medical And Biological Engineering And Computing*, 28(3):217–224, may 1990. doi:[10.1007/bf02442670](https://doi.org/10.1007/bf02442670). (cited on page 7)
- Miquel Calvo and Josep M. Oller. A Distance Between Multivariate Normal Distributions Based In An Embedding Into The Siegel Group. *Journal Of Multivariate Analysis*, 35(2):223–242, nov 1990. doi:[10.1016/0047-259x\(90\)90026-e](https://doi.org/10.1016/0047-259x(90)90026-e). (cited on page 3)
- Andrzej Cichocki, Sergio Cruces, and Shun-ichi Amari. Log-Determinant Divergences Revisited: Alpha-Beta And Gamma Log-Det Divergences. *Entropy*, 17(5):2988–3034, may 2015. doi:[10.3390/e17052988](https://doi.org/10.3390/e17052988). (cited on page 26)

-
- 540 D. R. Cox. The Regression Analysis Of Binary Sequences. *Journal Of The Royal Statistical Society Series B: Statistical Methodology*, 20(2):215–232, jul 1958. doi:10.1111/j.2517-6161.1958.tb00292.x.
541 (cited on page 9)
542
- 543 I. Csizár. A Class Of Measures Of Informativity Of Observation Channels. *Periodica Mathematica Hungarica*, 2(1-4):191–213, mar 1972. doi:10.1007/bf02018661. (cited on page 5)
544
- 545 Lieven De Lathauwer. An Extension Of The Generalized Svd For More Than Two Matrices. Technical
546 report, Katholieke Universiteit Leuven Group Science, Engineering and Technology nd E.E. Dept.
547 (ESAT) - SCD - SISTA,, Leuven, Belgium, 2011. (cited on page 7)
548
- 549 Gustavo Deco and Dragan Obradovic. *An Information-Theoretic Approach To Neural Computing*.
550 Springer New York, 1996. doi:10.1007/978-1-4612-4016-7. (cited on page 7)
551
- 552 James Weldon Demmel. Computing Stable Eigendecompositions Of Matrices. *Linear Algebra And
553 Its Applications*, 79:163–193, jul 1986. doi:10.1016/0024-3795(86)90298-3. (cited on page 6)
554
- 555 James Weldon Demmel and Bo Kågström. Computing Stable Eigendecompositions Of Matrix Pencils.
556 *Linear Algebra And Its Applications*, 88-89:139–186, apr 1987. doi:10.1016/0024-3795(87)90108-
557 x. (cited on page 2)
- 558 Yuxin Dong, Tieliang Gong, Hong Chen, Shuangyong Song, Weizhan Zhang, and Chen Li. How
559 Does Distribution Matching Help Domain Generalization: An Information-Theoretic Analysis. jun
560 2024. doi:10.48550/arXiv.2406.09745. (cited on page 7)
561
- 562 Harry Dym. *Linear Algebra In Action*. American Mathematical Society, 2023. doi:10.1090/gsm/232.
563 (cited on page 3)
- 564 Ky Fan. On A Theorem Of Weyl Concerning Eigenvalues Of Linear Transformations I. *Proceedings
565 Of The National Academy Of Sciences*, 35(11):652–655, nov 1949. doi:10.1073/pnas.35.11.652.
566 (cited on pages 2, 5, and 28)
567
- 568 Robert Geirhos, Jörn-Henrik Jacobsen, Claudio Michaelis, Richard Zemel, Wieland Brendel, Matthias
569 Bethge, and Felix A. Wichmann. Shortcut Learning In Deep Neural Networks. *Nature Machine
570 Intelligence*, 2:665–673, apr 2020. doi:10.1038/s42256-020-00257-z. (cited on page 8)
- 571 I. Gohberg, P. Lancaster, and L. Rodman. *Matrix Polynomials*. Society for Industrial and Applied
572 Mathematics, 2009. doi:10.1137/1.9780898719024. (cited on page 3)
573
- 574 Ian Goodfellow, Yoshua Bengio, and Aaron Courville. *Deep Learning*. MIT Press, 2017. (cited on
575 page 2)
- 576 Ishaan Gulrajani and David Lopez-Paz. In Search Of Lost Domain Generalization. *Arxiv*, 2020. (cited
577 on pages 1 and 7)
578
- 579 Nathan Halko, Per-Gunnar Martinsson, and Joel A. Tropp. Finding Structure With Randomness:
580 Probabilistic Algorithms For Constructing Approximate Matrix Decompositions. *Arxiv*, 53
581 (arXiv:909.4061):217–288, sep 2009. doi:10.48550/arXiv.0909.4061. (cited on page 9)
- 582 David Hilbert. *Grundzüge Einer Allgemeinen Theorie Der Linearen Integralgleichungen. (Zweite
583 Mitteilung)*, volume 1904, chapter Grundzüge einer allgemeinen Theorie der linearen Integral-
584 gleichungen, pp. 6–169. Springer Vienna, 1989. doi:10.1007/978-3-7091-9535-2_1. (cited on
585 page 4)
586
- 587 Geoffrey E. Hinton. Learning Multiple Layers Of Representation. *Trends In Cognitive Sciences*, 11
588 (10):428–434, oct 2007. doi:10.1016/j.tics.2007.09.004. (cited on page 9)
- 589 H. O. Hirschfeld. A Connection Between Correlation And Contingency. *Mathematical Proceedings Of
590 The Cambridge Philosophical Society*, 31(4):520–524, oct 1935. doi:10.1017/s0305004100013517.
591 (cited on page 7)
592
- 593 Reshad Hosseini and Suvrit Sra. Manifold Optimization For Gaussian Mixture Models. jun 2015.
doi:10.48550/arXiv.1506.07677. (cited on page 3)

594 Shao-Lun Huang, Anuran Makur, Gregory W. Wornell, and Lizhong Zheng. Universal Features For
595 High-Dimensional Learning And Inference. *Foundations And Trends In Communications And*
596 *Information Theory*, 21(1-2):1–299, 2024. doi:[10.1561/0100000107](https://doi.org/10.1561/0100000107). (cited on page 7)
597

598 Jörn-Henrik Jacobsen, Jens Behrmann, Richard Zemel, and Matthias Bethge. Excessive Invariance
599 Causes Adversarial Vulnerability. *Arxiv*, nov 2018. doi:[10.48550/arXiv.1811.00401](https://doi.org/10.48550/arXiv.1811.00401). (cited on
600 page 8)

601 Arthur Jacot, Franck Gabriel, and Clément Hongler. Neural Tangent Kernel: Convergence And
602 Generalization In Neural Networks. jun 2018. doi:[10.48550/arXiv.1806.07572](https://doi.org/10.48550/arXiv.1806.07572). (cited on page 7)
603

604 Tosio Kato. *Perturbation Theory For Linear Operators*. Springer Science & Business Media, 1995.
605 (cited on page 31)

606 Jivat Neet Kaur, Emre Kiciman, and Amit Sharma. Modeling The Data-Generating Process Is
607 Necessary For Out-Of-Distribution Generalization. jun 2022. doi:[10.48550/arXiv.2206.07837](https://doi.org/10.48550/arXiv.2206.07837).
608 (cited on page 7)
609

610 Liana Khamidullina, André L. F. de Almeida, and Martin Haardt. Multilinear Generalized Singular
611 Value Decomposition (MI-Gsvd) With Application To Coordinated Beamforming In Multi-User
612 MIMO Systems. In *Icassp 2020 - 2020 Ieee International Conference On Acoustics, Speech And*
613 *Signal Processing (Icassp)*, pp. 4587–4591, may 2020. doi:[10/gsvnbk](https://doi.org/10/gsvnbk). (cited on page 7)
614

615 Daniel Kressner. *Numerical Methods For General And Structured Eigenvalue Problems*. Springer-
616 Verlag, 2005. doi:[10.1007/3-540-28502-4](https://doi.org/10.1007/3-540-28502-4). (cited on pages 3 and 6)

617 David Krueger, Ethan Caballero, Joern-Henrik Jacobsen, Amy Zhang, Jonathan Binas, Dinghui
618 Zhang, Remi Le Priol, and Aaron Courville. Out-Of-Distribution Generalization Via Risk Extrapo-
619 lation (Rex). mar 2020. doi:[10.48550/arXiv.2003.00688](https://doi.org/10.48550/arXiv.2003.00688). (cited on pages 7 and 8)
620

621 Solomon Kullback. *Information Theory And Statistics*. Dover Publications, 1968. (cited on page 4)

622 Joseph-Louis Lagrange. *Mécanique Analytique*. Cambridge University Press, 2009.
623 doi:[10.1017/cbo9780511701788](https://doi.org/10.1017/cbo9780511701788). (cited on page 20)
624

625 H. O. Lancaster. The Structure Of Bivariate Distributions. *The Annals Of Mathematical Statistics*, 29
626 (3):719–736, sep 1958. doi:[10.1214/aoms/1177706532](https://doi.org/10.1214/aoms/1177706532). (cited on page 7)
627

628 Peter Lancaster and Leiba Rodman. Perturbation Theory For Discrete Algebraic Riccati Equations.
629 In *Algebraic Riccati Equations*, volume 1, pp. 329–332. Oxford University Press Oxford, sep 1995.
630 doi:[10.1093/oso/9780198537953.003.0014](https://doi.org/10.1093/oso/9780198537953.003.0014). (cited on page 6)

631 Kenneth Lange. *Mm Optimization Algorithms*. Society for Industrial and Applied Mathematics, 2016.
632 doi:[10.1137/1.9781611974409](https://doi.org/10.1137/1.9781611974409). (cited on page 4)
633

634 Yann LeCun. The Mnist Database Of Handwritten Digits. *Http://Yann. Lecun. Com/Exdb/Mnist/*,
635 1998. (cited on page 9)

636 Bo Li, Yezhen Wang, Shanghang Zhang, Dongsheng Li, Trevor Darrell, Kurt Keutzer, and Han Zhao.
637 Learning Invariant Representations And Risks For Semi-Supervised Domain Adaptation. oct 2020.
638 doi:[10.48550/arXiv.2010.04647](https://doi.org/10.48550/arXiv.2010.04647). (cited on page 7)
639

640 Yong Lin, Hanze Dong, Hao Wang, and Tong Zhang. Bayesian Invariant Risk Minimization. In *2022*
641 *Ieee/Cvf Conference On Computer Vision And Pattern Recognition (Cvpr)*, pp. 16000–16009, jun
642 2022. doi:[10.1109/cvpr52688.2022.01555](https://doi.org/10.1109/cvpr52688.2022.01555). (cited on page 7)

643 R. Linsker. Self-Organization In A Perceptual Network. *Computer*, 21(3):105–117, mar 1988.
644 doi:[10.1109/2.36](https://doi.org/10.1109/2.36). (cited on page 7)
645

646 Dong C. Liu and Jorge Nocedal. On The Limited Memory Bfgs Method For Large Scale Optimization.
647 *Mathematical Programming*, 45(1-3):503–528, aug 1989. doi:[10.1007/bf01589116](https://doi.org/10.1007/bf01589116). (cited on
page 9)

648 Divyat Mahajan, Shruti Tople, and Amit Sharma. Domain Generalization Using Causal Matching.
649 2020. doi:[10/gz3534](https://doi.org/10/gz3534). (cited on page 7)
650

651 Anuran Makur. *Information Contraction And Decomposition*. Thesis, Massachusetts Institute of
652 Technology, 2019. URL <https://dspace.mit.edu/handle/1721.1/122692>. (cited on page 7)
653

654 Per-Gunnar Martinsson, Vladimir Rokhlin, and Mark Tygert. A Randomized Algorithm For The
655 Decomposition Of Matrices. *Applied And Computational Harmonic Analysis*, 30(1):47–68, jan
656 2011. doi:[10.1016/j.acha.2010.02.003](https://doi.org/10.1016/j.acha.2010.02.003). (cited on page 9)

657 David F. Mayers, Gene H. Golub, and Charles F. van Loan. Matrix Computations. *Mathematics Of*
658 *Computation*, 47(175):376, jul 1986. doi:[10.2307/2008107](https://doi.org/10.2307/2008107). (cited on page 6)
659

660 Robb J. Muirhead. *Aspects Of Multivariate Statistical Theory*. Wiley, 1982.
661 doi:[10.1002/9780470316559](https://doi.org/10.1002/9780470316559). (cited on page 4)

662 C. C. Paige and M. A. Saunders. Towards A Generalized Singular Value Decomposition. *Siam*
663 *Journal On Numerical Analysis*, 18(3):398–405, jun 1981. doi:[10/bv55c5](https://doi.org/10/bv55c5). (cited on page 7)
664

665 Yury Polyanskiy and Yihong Wu. Dissipation Of Information In Channels With Input Constraints.
666 *Ieee Transactions On Information Theory*, 62(1):35–55, jan 2016. doi:[10.1109/tit.2015.2482978](https://doi.org/10.1109/tit.2015.2482978).
667 (cited on pages 7 and 41)

668 Yury Polyanskiy and Yihong Wu. *Information Theory*. Cambridge University Press, 2024.
669 doi:[10.1017/9781108966351](https://doi.org/10.1017/9781108966351). (cited on pages 4 and 24)
670

671 Sri Priya Ponnappalli, Michael A. Saunders, Charles F. Van Loan, and Orly Alter. A Higher-Order
672 Generalized Singular Value Decomposition For Comparison Of Global Mrna Expression From
673 Multiple Organisms. *Plos One*, 6(12):e28072, dec 2011. doi:[10.1371/journal.pone.0028072](https://doi.org/10.1371/journal.pone.0028072). (cited
674 on page 7)

675 Maxim Raginsky. Strong Data Processing Inequalities And Φ -Sobolev Inequalities For Discrete
676 Channels. nov 2014. doi:[10.48550/arXiv.1411.3575](https://doi.org/10.48550/arXiv.1411.3575). (cited on page 7)
677

678 Elan Rosenfeld, Pradeep Ravikumar, and Andrej Risteski. The Risks Of Invariant Risk Minimization.
679 oct 2020. doi:[10.48550/arXiv.2010.05761](https://doi.org/10.48550/arXiv.2010.05761). (cited on pages 2, 7, 8, and 33)

680 Sam Roweis and Carlos Brody. Linear Heteroencoders. *Gatsby Computational Neuroscience Unit,*
681 *Alexandra House: London, Uk*, 1999. (cited on page 7)
682

683 Shiori Sagawa, Pang Wei Koh, Tatsunori B. Hashimoto, and Percy Liang. Distributionally Ro-
684 bust Neural Networks For Group Shifts: On The Importance Of Regularization For Worst-Case
685 Generalization. nov 2019. doi:[10.48550/arXiv.1911.08731](https://doi.org/10.48550/arXiv.1911.08731). (cited on pages 2 and 9)

686 Erhard Schmidt. Zur Theorie Der Linearen Und Nichtlinearen Integralgleichungen. *Mathematische*
687 *Annalen*, 63(4):433–476, dec 1907. doi:[10.1007/bf01449770](https://doi.org/10.1007/bf01449770). (cited on page 7)
688

689 Harshay Shah, Kaustav Tamuly, Aditi Raghunathan, Prateek Jain, and Praneeth Netrapalli. The
690 Pitfalls Of Simplicity Bias In Neural Networks. jun 2020. doi:[10.48550/arXiv.2006.07710](https://doi.org/10.48550/arXiv.2006.07710). (cited
691 on page 9)

692 C. E. Shannon. A Mathematical Theory Of Communication. *Bell System Technical Journal*, 27(3):
693 379–423, jul 1948. doi:[10.1002/j.1538-7305.1948.tb01338.x](https://doi.org/10.1002/j.1538-7305.1948.tb01338.x). (cited on page 4)
694

695 Carl Ludwig Siegel. Symplectic Geometry. *American Journal Of Mathematics*, 65(1):1, jan 1943.
696 doi:[10.2307/2371774](https://doi.org/10.2307/2371774). (cited on page 3)

697 Sanjeet Singh, Pankaj Gupta 0001, and Davinder Bhatia. Multiparametric Sensitivity Analysis Of The
698 Constraint Matrix In Linear-Plus-Linear Fractional Programming Problem. *Appl. Math. Comput.*,
699 170(2):1243–1260, 2005. doi:[10.1016/J.AMC.2005.01.016](https://doi.org/10.1016/J.AMC.2005.01.016). (cited on page 3)
700

701 G. W. Stewart, J. W. Stewart, and Ji-Guang Sun. *Matrix Perturbation Theory*. Elsevier Science, 1990.
(cited on page 6)

702 Yangfeng Su and Zhaojun Bai. Solving Rational Eigenvalue Problems Via Linearization. *Siam*
703 *Journal On Matrix Analysis And Applications*, 32(1):201–216, jan 2011. doi:[10.1137/090777542](https://doi.org/10.1137/090777542).
704 (cited on page 4)
705

706 Ji-guang Sun. Eigenvalues Of Rayleigh Quotient Matrices. *Numerische Mathematik*, 59(1):603–614,
707 dec 1991. doi:[10.1007/bf01385798](https://doi.org/10.1007/bf01385798). (cited on page 4)

708 Remi Tachet, Han Zhao, Yu-Xiang Wang, and Geoff Gordon. Domain Adaptation With Conditional
709 Distribution Matching And Generalized Label Shift. mar 2020. doi:[10.48550/arXiv.2003.04475](https://doi.org/10.48550/arXiv.2003.04475).
710 (cited on page 7)
711

712 Charles F. Van Loan. Generalizing The Singular Value Decomposition. *Siam Journal On Numerical*
713 *Analysis*, 13(1):76–83, mar 1976. doi:[10/cmmbjs](https://doi.org/10/cmmbjs). (cited on page 7)

714 Zihao Wang and Victor Veitch. A Unified Causal View Of Domain Invariant Representation Learning.
715 In *Icml: Workshop On Spurious Correlations, Invariance And Stability*, 2022. (cited on page 7)
716

717 Yuanzhe Xi and Y. Saad. Least-Squares Rational Filters For The Solution Of Interior Eigenvalue
718 Problems. 2015. (cited on page 4)

719 Greg Yang. Tensor Programs I: Wide Feedforward Or Recurrent Neural Networks Of Any Architecture
720 Are Gaussian Processes. oct 2019. doi:[10.48550/arXiv.1910.12478](https://doi.org/10.48550/arXiv.1910.12478). (cited on page 7)
721

722 Haotian Ye, Chuanlong Xie, Tianle Cai, Ruichen Li, Zhenguo Li, and Liwei Wang. Towards A Theoret-
723 ical Framework Of Out-Of-Distribution Generalization. jun 2021. doi:[10.48550/arXiv.2106.04496](https://doi.org/10.48550/arXiv.2106.04496).
724 (cited on page 7)

725 Matthew D Zeiler and Rob Fergus. Visualizing And Understanding Convolutional Networks. nov
726 2013. doi:[10.48550/arXiv.1311.2901](https://doi.org/10.48550/arXiv.1311.2901). (cited on page 6)
727

728 Sanfeng Zhang, Xinyi Liu, Zihao Qi, Xingchen Yan, and Wang Yang. Gi-Graph: A Generative
729 Invariant Graph Learning Scheme Towards Out-Of-Distribution Generalization. *Ieee Transactions*
730 *On Knowledge And Data Engineering*, 35:1–15, 2025. doi:[10.1109/tkde.2025.3592640](https://doi.org/10.1109/tkde.2025.3592640). (cited on
731 pages 7 and 8)

732 Xiao Zhou, Yong Lin, Weizhong Zhang, and Tong Zhang. Sparse Invariant Risk Minimization. In
733 *Proceedings Of The 39Th International Conference On Machine Learning*, pp. 27222–27244, 2022.
734 (cited on page 7)

735 Yi Zhou and Yingbin Liang. Critical Points Of Neural Networks: Analytical Forms And Landscape
736 Properties. oct 2017. doi:[10.48550/arXiv.1710.11205](https://doi.org/10.48550/arXiv.1710.11205). (cited on page 7)
737
738
739
740
741
742
743
744
745
746
747
748
749
750
751
752
753
754
755

756
757
758
759
760
761
762
763
764
765
766
767
768
769
770
771
772
773
774
775
776
777
778
779
780
781
782
783
784
785
786
787
788
789
790
791
792
793
794
795
796
797
798
799
800
801
802
803
804
805
806
807
808
809

Appendices

Contents

A Two states model example	15
B Parametrized spectral ratio example	16
C Constructing Q	18
C.1 Phi-map is injective	18
C.2 The image of the Phi-map is a subset of the pd cone	18
C.3 Q is convex in the pd cone	18
C.4 Algebraic structure of Q	19
C.5 Trace and determinants of ratio of Phi-embedding	19
D Proofs	20
D.1 Theorem 1	20
D.2 Lemma 1	20
D.3 Proposition 1	21
D.4 Theorem 2	25
D.5 Corollary 1	28
D.6 Theorem 3	28
D.7 Proposition 2	32
E Mathematical elements	37
E.1 Schur complement	37
E.2 LDU and UDL Decomposition	37
E.3 Simultaneous diagonalization of SPD matrices by congruence	39
E.4 Spectral equivalences	39
E.5 Concavity of the log-determinant	40
E.6 Necessary and sufficient solutions for a Pseudo-Inverse solution	41
E.7 f-divergence	41

A TWO STATES MODEL EXAMPLE

Consider a state space with two elements $\mathcal{L} = \{0, 1\}$. Let $\mathbf{I} \sim \text{Ber}(p)$. When $\mathbf{I} = 1$, $\mathbf{x} \sim \mathcal{N}\left(\begin{bmatrix} 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}\right)$. When $\mathbf{I} = 2$, $\mathbf{x} \sim \mathcal{N}\left(\begin{bmatrix} \mu \\ 0 \end{bmatrix}, \begin{bmatrix} 1 & \rho \\ \rho & 1 \end{bmatrix}\right)$, see figure 3. Using definition 1 with \mathbf{H} equal to the covariance and \mathbf{G} the mean vector, we see that the Φ -embedding of the input stream in the first state is simply the identity matrix, $\Phi_1 = \mathbf{I}_3$, that of the second state is given by,

$$\Phi_2 = \begin{bmatrix} 1 + \mu^2 & \rho & \mu \\ \rho & 1 & 0 \\ \mu & 0 & 1 \end{bmatrix}.$$

The expected state is given by,

$$\begin{aligned} \mathbb{E}[\Phi_1] &= \mathbb{P}(\mathbf{I} = 1)\Phi_1 + \mathbb{P}(\mathbf{I} = 2)\Phi_2, \\ &= p\Phi_1 + (1 - p)\Phi_2, \\ &= \begin{bmatrix} 1 + \mu^2(1 - p) & \rho(1 - p) & \mu(1 - p) \\ \rho(1 - p) & 1 & 0 \\ \mu(1 - p) & 0 & 1 \end{bmatrix}. \end{aligned}$$

810 Let us now form the expected ratio,
811

$$812 \quad \mathbb{E}[\mathbb{E}[\Phi_1] : \Phi_1] = p (\mathbb{E}[\Phi_1] : \mathbf{I}) + (1-p) (\mathbb{E}[\Phi_1] : \Phi_2),$$

$$813 \quad = p \mathbb{E}[\Phi_1] + (1-p) (\mathbb{E}[\Phi_1] : \Phi_2).$$

814 We need to evaluate the ratio $\mathbb{E}[\Phi_1] : \Phi_2$. We have,
815

$$816 \quad \Phi_2 = \mathbf{U} \left(\begin{bmatrix} \mu \\ 0 \end{bmatrix} \right) \left(\begin{bmatrix} 1 & \rho \\ \rho & 1 \end{bmatrix} \oplus 1 \right) \mathbf{U} \left(\begin{bmatrix} \mu \\ 0 \end{bmatrix} \right)^\top,$$

$$817 \quad = \begin{bmatrix} 1 & 0 & \mu \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} 1 & \rho & 0 \\ \rho & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ \mu & 0 & 1 \end{bmatrix}$$

818 Now by the spectral theorem,
819

$$820 \quad \begin{bmatrix} 1 & \rho \\ \rho & 1 \end{bmatrix} = \mathbf{V} \mathbf{\Lambda} \mathbf{V}^\top,$$

$$821 \quad = \begin{bmatrix} -1 & 1 \\ 1 & 1 \end{bmatrix} \begin{bmatrix} 1-\rho & 0 \\ 0 & 1+\rho \end{bmatrix} \begin{bmatrix} -1 & 1 \\ 1 & 1 \end{bmatrix}$$

822 Hence,
823

$$824 \quad \begin{bmatrix} 1 & \rho \\ \rho & 1 \end{bmatrix}^{-\frac{1}{2}} = \begin{bmatrix} -1 & 1 \\ 1 & 1 \end{bmatrix} \begin{bmatrix} \frac{1}{\sqrt{1-\rho}} & 0 \\ 0 & \frac{1}{\sqrt{1+\rho}} \end{bmatrix} \begin{bmatrix} -1 & 1 \\ 1 & 1 \end{bmatrix}$$

825 Using the definition of the ratio theorem 1 we compute,
826

$$827 \quad \mathbb{E}[\Phi_1] : \Phi_2 = \begin{bmatrix} 1 - \frac{p(\mu^2\sqrt{1-\rho^2} + \mu^2 + 2\rho^2)}{2(\rho^2-1)} & \frac{(\mu^2+2)p\rho}{2(\rho^2-1)} & -\frac{\mu p(\sqrt{1-\rho} + \sqrt{\rho+1})}{2\sqrt{1-\rho^2}} \\ \frac{(\mu^2+2)p\rho}{2(\rho^2-1)} & \frac{\mu^2 p \sqrt{1-\rho^2} + \mu^2(-\rho) - 2p\rho^2 + 2\rho^2 - 2}{2(\rho^2-1)} & -\frac{\mu p(\sqrt{1-\rho} - \sqrt{\rho+1})}{2\sqrt{1-\rho^2}} \\ -\frac{\mu p(\sqrt{1-\rho} + \sqrt{\rho+1})}{2\sqrt{1-\rho}\sqrt{\rho+1}} & -\frac{\mu p(\sqrt{1-\rho} - \sqrt{\rho+1})}{2\sqrt{1-\rho}\sqrt{\rho+1}} & 1 \end{bmatrix}.$$

828 We can then compute the dependence on the state \mathbf{I} . By theorem 2
829

$$830 \quad S(\mathbf{x}; \mathbf{I}) = \frac{1}{2} \text{Tr} [\mathbb{E}[\Phi_1] : \Phi_1],$$

$$831 \quad = \frac{p(1-p)(\mu^2\rho^2 - 2\mu^2 - 2\rho^2)}{2(\rho^2-1)}.$$

832 B PARAMETRIZED SPECTRAL RATIO EXAMPLE

833 Consider two Gaussian distributions, $\mathbb{P}_1 = \mathcal{N}\left(\begin{bmatrix} \frac{1}{\sqrt{2+\delta}} \\ \frac{\delta}{2} \end{bmatrix}, \begin{bmatrix} 1-\delta^2 & 0 \\ 0 & 1 \end{bmatrix}\right)$, $\mathbb{P}_2 = \mathcal{N}\left(\begin{bmatrix} \frac{1}{\sqrt{2-\delta}} \\ -\frac{\delta}{2} \end{bmatrix}, \begin{bmatrix} 1 & 0 \\ 0 & 1+\delta^2 \end{bmatrix}\right)$.
834 These distributions exist as long as $\delta \in (-1, 1)$. By lemma 1. Forming their respective canonical
835 embeddings. We have, $\lambda(\Phi_1 : \Phi_2) = \left\{ 1 - \delta^2, 1 - \frac{|\delta|}{\sqrt{\delta^2+1}}, 1 + \frac{|\delta|}{\sqrt{\delta^2+1}} \right\}$. Figure 6 shows the evolution
836 of the spectrum of the ratio as a function of δ .
837
838
839
840
841
842
843
844
845
846
847
848

864
865
866
867
868
869
870
871
872
873
874
875
876
877
878
879
880
881
882
883
884
885
886
887
888
889
890
891
892
893
894
895
896
897
898
899
900
901
902
903
904
905
906
907
908
909
910
911
912
913
914
915
916
917

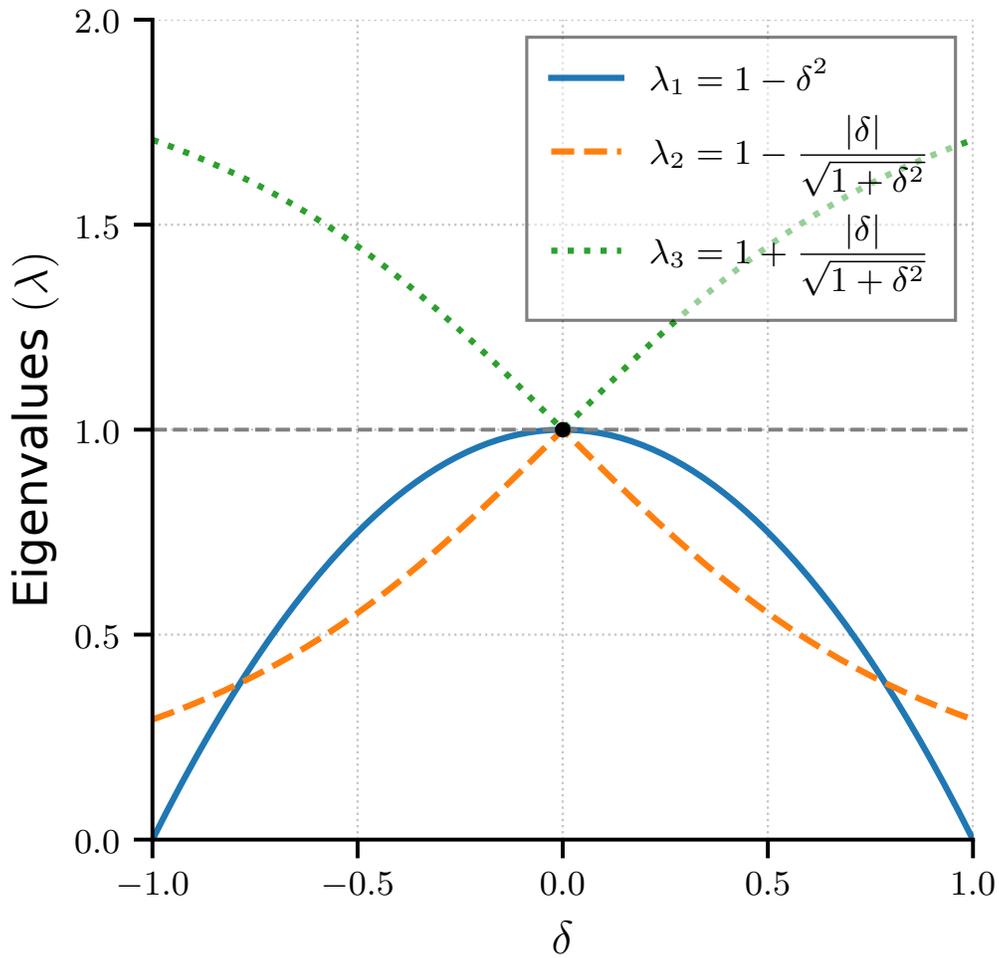


Figure 6: Spectrum of the parametrized ratio of Φ -embeddings.

918
919
920
921
922
923
924
925
926
927
928
929
930
931
932
933
934
935
936
937
938
939
940
941
942
943
944
945
946
947
948
949
950
951
952
953
954
955
956
957
958
959
960
961
962
963
964
965
966
967
968
969
970
971

C CONSTRUCTING Q

C.1 PHI-MAP IS INJECTIVE

Proposition 3. *The Φ -map is injective.*

Proof. Indeed,

$$\begin{aligned} \begin{bmatrix} \mathbf{H}_1 + \mathbf{G}_1 \mathbf{G}_1^\top & \mathbf{G}_1 \\ \mathbf{G}_1^\top & \mathbf{I}_n \end{bmatrix} = \begin{bmatrix} \mathbf{H}_2 + \mathbf{G}_2 \mathbf{G}_2^\top & \mathbf{G}_2 \\ \mathbf{G}_2^\top & \mathbf{I}_n \end{bmatrix} &\iff \\ \begin{bmatrix} \mathbf{H}_1 - \mathbf{H}_2 + \mathbf{G}_1 \mathbf{G}_1^\top - \mathbf{G}_2 \mathbf{G}_2^\top & \mathbf{G}_1 - \mathbf{G}_2 \\ \mathbf{G}_1^\top - \mathbf{G}_2^\top & \mathbf{0} \end{bmatrix} = \begin{bmatrix} \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix} &\iff, \\ (\mathbf{G}_1, \mathbf{H}_1) = (\mathbf{G}_2, \mathbf{H}_2). & \end{aligned}$$

□

Hence, the Φ -map defines an embedding, $\mathcal{S}_\geq(m) \times \mathcal{M}_{m \times n}(\mathbb{R}) \hookrightarrow \mathcal{S}_\geq(m+n)$.

Next we define the Φ -set.

Definition 3 (\mathcal{Q} -set). \mathcal{Q} is the set $\mathcal{Q}(m, n) := \{\Phi(\mathbf{H}, \mathbf{G}) : \mathbf{H} \in \mathcal{S}_>(m), \mathbf{G} \in \mathcal{M}_{m \times n}(\mathbb{R})\}$.

C.2 THE IMAGE OF THE PHI-MAP IS A SUBSET OF THE PD CONE

The next proposition show that the \mathcal{Q} is indeed a subset of some cone of symmetric positive definite matrices.

Proposition 4. $\mathcal{Q}(m, n) \subset \mathcal{S}_>(m+n)$.

Proof. $\Phi := \Phi(\mathbf{H}, \mathbf{G}) > \mathbf{0} \iff \mathbf{I}_n > \mathbf{0}$ and $\Phi/\mathbf{I}_n = \mathbf{H} > \mathbf{0}$.

□

The next proposition shows that \mathcal{Q} is convex.

C.3 Q IS CONVEX IN THE PD CONE

Proposition 5. *Let $\pi \in \Delta^{k-1}$, for all $i \in [k]$, $\mathbf{H}_i, \mathbf{G}_i \in \mathcal{S}_>(m) \times \mathcal{M}_{m \times n}(\mathbb{R})$. Define, $\bar{\mathbf{G}} = \sum_{i=1}^k \pi_i \mathbf{G}_i$, $\bar{\mathbf{H}} = \sum_{i=1}^k \pi_i (\mathbf{H}_i + (\mathbf{G}_i - \bar{\mathbf{G}})(\mathbf{G}_i - \bar{\mathbf{G}})^\top)$. Then, $\Phi(\bar{\mathbf{H}}, \bar{\mathbf{G}}) = \sum_{i=1}^k \pi_i \Phi(\mathbf{H}_i, \mathbf{G}_i)$.*

Proof. First, since $\forall i \in [k], (\mathbf{G}_i - \bar{\mathbf{G}})(\mathbf{G}_i - \bar{\mathbf{G}})^\top \geq \mathbf{0}$ and $\mathbf{H}_i > \mathbf{0}$, we have that $\forall i \in [k], \pi_i (\mathbf{H}_i + (\mathbf{G}_i - \bar{\mathbf{G}})(\mathbf{G}_i - \bar{\mathbf{G}})^\top) > \mathbf{0}$ and hence, $\bar{\mathbf{H}} > \mathbf{0}$. Hence, $\Phi(\bar{\mathbf{H}}, \bar{\mathbf{G}}) \in \mathcal{Q}(m, n)$. Expanding Φ ,

$$\begin{aligned} \Phi(\bar{\mathbf{H}}, \bar{\mathbf{G}}) &= \begin{bmatrix} \bar{\mathbf{H}} + \bar{\mathbf{G}} \bar{\mathbf{G}}^\top & \bar{\mathbf{G}} \\ \bar{\mathbf{G}}^\top & \mathbf{I} \end{bmatrix}, \\ &= \begin{bmatrix} \sum_{i=1}^k \pi_i (\mathbf{H}_i + (\mathbf{G}_i - \bar{\mathbf{G}})(\mathbf{G}_i - \bar{\mathbf{G}})^\top) + \bar{\mathbf{G}} \bar{\mathbf{G}}^\top & \bar{\mathbf{G}} \\ \bar{\mathbf{G}}^\top & \mathbf{I} \end{bmatrix}, \\ &= \begin{bmatrix} \sum_{i=1}^k \pi_i (\mathbf{H}_i + \mathbf{G}_i \mathbf{G}_i^\top) - \bar{\mathbf{G}} \bar{\mathbf{G}}^\top + \bar{\mathbf{G}} \bar{\mathbf{G}}^\top & \bar{\mathbf{G}} \\ \bar{\mathbf{G}}^\top & \mathbf{I} \end{bmatrix}, \\ &= \begin{bmatrix} \sum_{i=1}^k \pi_i (\mathbf{H}_i + \mathbf{G}_i \mathbf{G}_i^\top) & \bar{\mathbf{G}} \\ \bar{\mathbf{G}}^\top & \mathbf{I} \end{bmatrix}, \\ &= \sum_{i=1}^k \pi_i \begin{bmatrix} \mathbf{H}_i + \mathbf{G}_i \mathbf{G}_i^\top & \mathbf{G}_i \\ \mathbf{G}_i^\top & \mathbf{I} \end{bmatrix}, \\ &= \sum_{i=1}^k \pi_i \Phi(\mathbf{H}_i, \mathbf{G}_i). \end{aligned}$$

□

972
973
974
975
976
977
978
979
980
981
982
983
984
985
986
987
988
989
990
991
992
993
994
995
996
997
998
999
1000
1001
1002
1003
1004
1005
1006
1007
1008
1009
1010
1011
1012
1013
1014
1015
1016
1017
1018
1019
1020
1021
1022
1023
1024
1025

Corollary 3. $\mathcal{Q}(m, n)$ is a convex subset of $\mathcal{S}_{>}(m + n)$

Proof. Follows from proposition 5. □

C.4 ALGEBRAIC STRUCTURE OF \mathcal{Q}

The next lemma expresses the ratio of two element in \mathcal{Q} .

Lemma 2. Let $(\mathbf{H}_1, \mathbf{G}_1), (\mathbf{H}_2, \mathbf{G}_2) \in \mathcal{S}_{>}(m) \times \mathcal{M}_{m \times n}(\mathbb{R})$, and define, $\Phi_1 := \Phi(\mathbf{H}_1, \mathbf{G}_1)$, $\Phi_2 := \Phi(\mathbf{H}_2, \mathbf{G}_2)$. Then, $\Phi_1 : \Phi_2 = \Phi(\mathbf{H}_2^{-\frac{1}{2}} \mathbf{H}_1 \mathbf{H}_2^{-\frac{1}{2}}, \mathbf{H}_2^{-\frac{1}{2}} (\mathbf{G}_1 - \mathbf{G}_2))$ and, $\Phi_2 : \Phi_1 = \Phi(\mathbf{H}_1^{-\frac{1}{2}} \mathbf{H}_2 \mathbf{H}_1^{-\frac{1}{2}}, \mathbf{H}_1^{-\frac{1}{2}} (\mathbf{G}_2 - \mathbf{G}_1))$.

Proof. By definition of the $:$,

$$\begin{aligned} \Phi_1 : \Phi_2 &= \begin{bmatrix} \mathbf{H}_2^{-\frac{1}{2}} (\mathbf{G}_1 \mathbf{G}_1^\top + \mathbf{G}_2 \mathbf{G}_2^\top + \mathbf{H}_1 - \mathbf{G}_2 \mathbf{G}_1^\top - \mathbf{G}_1 \mathbf{G}_2^\top) \mathbf{H}_2^{-\frac{1}{2}} & \mathbf{H}_2^{-\frac{1}{2}} (\mathbf{G}_1 - \mathbf{G}_2) \\ (\mathbf{G}_1^\top - \mathbf{G}_2^\top) \mathbf{H}_2^{-\frac{1}{2}} & \mathbf{I}_n \end{bmatrix}, \\ &= \begin{bmatrix} \mathbf{H}_2^{-\frac{1}{2}} (\mathbf{H}_1 + (\mathbf{G}_1 - \mathbf{G}_2)(\mathbf{G}_1 - \mathbf{G}_2)^\top) \mathbf{H}_2^{-\frac{1}{2}} & \mathbf{H}_2^{-\frac{1}{2}} (\mathbf{G}_1 - \mathbf{G}_2) \\ (\mathbf{G}_1 - \mathbf{G}_2)^\top \mathbf{H}_2^{-\frac{1}{2}} & \mathbf{I}_n \end{bmatrix}, \\ &= \Phi(\mathbf{H}_2^{-\frac{1}{2}} \mathbf{H}_1 \mathbf{H}_2^{-\frac{1}{2}}, \mathbf{H}_2^{-\frac{1}{2}} (\mathbf{G}_1 - \mathbf{G}_2)). \end{aligned}$$

The expression for $\Phi_2 : \Phi_1$ follows by symmetry. □

Note that, akin to division on positive reals, $:$ is not commutative.

We now start to endowing \mathcal{Q} with algebraic structure. First recall,

Definition 4 (Magma). A magma, $(\mathcal{Q}, :)$ is a set \mathcal{Q} closed under a binary operation $:$.

We show that \mathcal{Q} is a magma,

Corollary 4 (\mathcal{Q} is a magma.). $\mathcal{Q}(m, n)$ is closed under $:$.

Proof. Let $\Phi_1, \Phi_2 \in \mathcal{Q}(m, n)$. By proposition 3 there exists, $(\mathbf{H}_1, \mathbf{G}_1), (\mathbf{H}_2, \mathbf{G}_2) \in \mathcal{S}_{>}(m) \times \mathcal{M}_{m \times n}(\mathbb{R})$, such that $\Phi_1 = \Phi(\mathbf{H}_1, \mathbf{G}_1)$, $\Phi_2 = \Phi(\mathbf{H}_2, \mathbf{G}_2)$. The conclusion follows from Lemma 2. □

C.5 TRACE AND DETERMINANTS OF RATIO OF PHI-EMBEDDING

Let's express the determinant and trace of ratios of Φ -embeddings.

Proposition 6. Let, $\Phi_i := \Phi(\mathbf{H}_i, \mathbf{G}_i)$, for $i \in \{1, 2\}$. The

- (i) $\det[\Phi_1 : \Phi_2] = \det[\mathbf{H}_2^{-1} \mathbf{H}_1]$,
- (ii) $\text{Tr}[\Phi_1 : \Phi_2] = \text{Tr}[\mathbf{H}_2^{-1} \mathbf{H}_1] + \text{Tr}[(\mathbf{G}_2 - \mathbf{G}_1)^\top \mathbf{H}_2^{-1} (\mathbf{G}_2 - \mathbf{G}_1)] + n$.

Proof. (i) By proposition 26 $\lambda(\Phi_1 : \Phi_2) = \lambda(\Phi_1, \Phi_2)$. Moreover, $\lambda(\Phi_1, \Phi_2) = \lambda(\Phi_2^{-1} \Phi_1)$. Since the determinant is equal to the product of the eigenvalues we have, $\det[\Phi_1 : \Phi_2] = \det[\Phi_2^{-1} \Phi_1] = \det[\Phi_2^{-1}] \det[\Phi_1] = \det[\Phi_2]^{-1} \det[\Phi_1]$, but $\det[\Phi_i] = \det[\mathbf{H}_i]$, hence $\det[\Phi_1 : \Phi_2] = \det[\mathbf{H}_2^{-1} \mathbf{H}_1]$. (ii) Recalling that the trace of block diagonal matrix is equal to sum of the traces of the blocks and that the second block is the identity and its trace is equal to n . It is enough to compute the upper left hand block of $\Phi_1 : \Phi_2$. By lemma 2,

$$\begin{aligned} (\Phi_1 : \Phi_2)[: n, : n] &= \mathbf{H}_2^{-\frac{1}{2}} \mathbf{H}_1 \mathbf{H}_2^{-\frac{1}{2}} + \mathbf{H}_2^{-\frac{1}{2}} \mathbf{G}_1 \mathbf{G}_1^\top \mathbf{H}_2^{-\frac{1}{2}} + \mathbf{H}_2^{-\frac{1}{2}} \mathbf{G}_2 \mathbf{G}_2^\top \mathbf{H}_2^{-\frac{1}{2}} \\ &\quad - \mathbf{H}_2^{-\frac{1}{2}} \mathbf{G}_2 \mathbf{G}_1^\top \mathbf{H}_2^{-\frac{1}{2}} - \mathbf{H}_2^{-\frac{1}{2}} \mathbf{G}_1 \mathbf{G}_2^\top \mathbf{H}_2^{-\frac{1}{2}} \\ &= \mathbf{H}_2^{-\frac{1}{2}} (\mathbf{H}_1 + (\mathbf{G}_2 - \mathbf{G}_1)(\mathbf{G}_2 - \mathbf{G}_1)^\top) \mathbf{H}_2^{-\frac{1}{2}}. \end{aligned}$$

Now, by proposition 26, $\text{Tr}[\Phi_1 : \Phi_2] - \text{Tr}[\mathbf{I}_n] = \text{Tr}[\mathbf{H}_2^{-\frac{1}{2}} \mathbf{H}_1 \mathbf{H}_2^{-\frac{1}{2}}] + \text{Tr}[\mathbf{H}_2^{-1}(\mathbf{G}_2 - \mathbf{G}_1)(\mathbf{G}_2 - \mathbf{G}_1)^\top] = \text{Tr}[\mathbf{H}_2^{-1} \mathbf{H}_1] + \text{Tr}[(\mathbf{G}_2 - \mathbf{G}_1)^\top \mathbf{H}_2^{-1}(\mathbf{G}_2 - \mathbf{G}_1)]$. \square

D PROOFS

D.1 THEOREM 1

Theorem 1 (Operator division algebra). *Let $(\mathbf{H}_1, \mathbf{G}_1) \in \mathcal{S}_>(m) \times \mathcal{M}_{m \times n}(\mathbb{R})$, and $\Phi_1 := \Phi(\mathbf{H}_1, \mathbf{G}_1)$. Then $\mathcal{Q}(m, n)$ is convex and closed under the binary ratio operations:*

$(\Phi_1, \Phi_{1'}) \mapsto \Phi_1 : \Phi_{1'} := \left(\mathbf{H}_{1'}^{-\frac{1}{2}} \oplus \mathbf{I}_n \right) \mathbf{U}(-\mathbf{G}_{1'}) \Phi_1 \mathbf{U}(-\mathbf{G}_{1'})^\top \left(\mathbf{H}_{1'}^{-\frac{1}{2}} \oplus \mathbf{I}_n \right)$. Moreover,

(i) $\Phi_1 : \mathbf{I} = \Phi_1$,

(ii) $\Phi_1 : \Phi_{1'} = \Phi(\mathbf{Q}_{1'}, \Delta_{1'})$, with $\mathbf{Q}_{1'} := \mathbf{H}_{1'}^{-\frac{1}{2}} \mathbf{H}_1 \mathbf{H}_{1'}^{-\frac{1}{2}}$ and $\Delta_{1'} := \mathbf{H}_{1'}^{-\frac{1}{2}}(\mathbf{G}_1 - \mathbf{G}_{1'})$,

(iii) $\mathbb{E}[\Phi_1 : \Phi_{1'}] = \Phi(\mathbb{E}[\mathbf{Q}_{1'}] + \mathbb{V}[\Delta_{1'}], \mathbb{E}[\Delta_{1'}])$.

Proof. The convexity follows from corollary 3. Closure under $:$ follows from corollary 4. (i) and (ii) follow from lemma 2. (iii) follows from (ii) and proposition 5. \square

D.2 LEMMA 1

Lemma 1. *Let $\mathbf{H} \in \mathcal{S}_>(m)$, $\mathbf{G} \in \mathcal{M}_{m \times n}(\mathbb{R})$, $\Phi := \Phi(\mathbf{H}, \mathbf{G}) \in \mathcal{Q}(m, n)$, and $\mathbf{w}^\top = [\mathbf{w}_1^\top \quad \mathbf{w}_2^\top]$ be a unit normalized eigenvector of Φ . Then any eigenpair (λ, \mathbf{W}) of Φ verifies,*

(i) *If $\lambda = 1$, $\mathbf{G}^\top \mathbf{w}_1 = 0$ and $(\mathbf{H} - \mathbf{I})\mathbf{w}_1 = -\mathbf{G}\mathbf{w}_2$.*

(ii) *If $\lambda \neq 1$, (λ, \mathbf{w}) is an eigenpair of Φ if and only if (λ, \mathbf{w}_1) is an eigenpair of the rational eigenvalue problem, $\mathbf{H}\mathbf{w}_1 = \lambda \left(\mathbf{I} + \frac{1}{1-\lambda} \mathbf{G}\mathbf{G}^\top \right) \mathbf{w}_1$, and $\mathbf{G}^\top \mathbf{w}_1 = (1-\lambda)\mathbf{w}_2$.*

Proof. Form the Rayleigh ratio, $\mathbf{w} \mapsto R[\mathbf{w}] := \frac{\mathbf{w}^\top \Phi \mathbf{w}}{\mathbf{w}^\top \mathbf{w}}$. The critical points(values) of R are the eigenvector (eigenvalues) of Φ . Partition, $\mathbf{w} = \begin{bmatrix} \mathbf{w}_1 \\ \mathbf{w}_2 \end{bmatrix}$. We'll proceed by block ascent over the eigenblocks \mathbf{w}_2 and \mathbf{w}_1 . Recall, $\Phi = \begin{bmatrix} \mathbf{H} + \mathbf{G}\mathbf{G}^\top & \mathbf{G} \\ \mathbf{G}^\top & \mathbf{I}_n \end{bmatrix}$. Maximization of R over $\{\mathbf{w} \neq \mathbf{0} \mid \mathbf{w} \in \mathbb{R}^{m+n}\}$, is equivalent that of the following program,

$$\min_{(\mathbf{w}_1, \mathbf{w}_2) \in \mathbb{R}^m \times \mathbb{R}^n} -\mathbf{w}^\top \Phi \mathbf{w},$$

$$\text{Such that, } \mathbf{w}^\top \mathbf{w} = 1.$$

We have, $\mathbf{w}^\top \Phi \mathbf{w} = \mathbf{w}_1^\top \mathbf{H} \mathbf{w}_1 + \mathbf{w}_1^\top \mathbf{G} \mathbf{G}^\top \mathbf{w}_1 + \mathbf{w}_2^\top \mathbf{G}^\top \mathbf{w}_1 + \mathbf{w}_1^\top \mathbf{G} \mathbf{w}_2 + \mathbf{w}_2^\top \mathbf{w}_2 = \mathbf{w}_1^\top \mathbf{H} \mathbf{w}_1 + (\mathbf{w}_2 + \mathbf{G}^\top \mathbf{w}_1)^\top (\mathbf{w}_2 + \mathbf{G}^\top \mathbf{w}_1)$. Form the Lagrangian (Lagrange, 2009),

$$L[\mathbf{w}_1, \mathbf{w}_2; \lambda] = -\mathbf{w}_1^\top \mathbf{H} \mathbf{w}_1 - (\mathbf{w}_2 + \mathbf{G}^\top \mathbf{w}_1)^\top (\mathbf{w}_2 + \mathbf{G}^\top \mathbf{w}_1) + \lambda(\mathbf{w}_1^\top \mathbf{w}_1 + \mathbf{w}_2^\top \mathbf{w}_2 - 1). \quad (1)$$

We will proceed by block coordinate descent starting with \mathbf{w}_2 . Taking the first differential, $dL(\mathbf{w}_2; d\mathbf{w}_2) = 2 \text{Tr} [(-\mathbf{w}_1^\top \mathbf{G} + \mathbf{w}_2^\top) + \lambda \mathbf{w}_2^\top] d\mathbf{w}_2$. The first order condition holds if and only if, $\mathbf{w}_2 + \mathbf{G}^\top \mathbf{w}_1 = \lambda \mathbf{w}_2 \iff \mathbf{G}^\top \mathbf{w}_1 = (\lambda - 1)\mathbf{w}_2$. Taking first order differential of equation (1) with respect to \mathbf{w}_1 ,

$$dL(\mathbf{w}_1, d\mathbf{w}_1) = -2(\mathbf{w}_1^\top \mathbf{H} + \mathbf{w}_1^\top \mathbf{G} \mathbf{G}^\top + \mathbf{w}_2^\top \mathbf{G}^\top - \lambda \mathbf{w}_1^\top) d\mathbf{w}_1. \quad (2)$$

We will consider two separate cases: $\lambda = 1$ and $\lambda \neq 1$.

If $\lambda = 1$, then $\mathbf{G}^\top \mathbf{w}_1 = 0$. Hence, $dL(\mathbf{w}_1; d\mathbf{w}_1) = 0 \iff (\mathbf{H} - \mathbf{I})\mathbf{w}_1 = -\mathbf{G}\mathbf{w}_2$.

If $\lambda \neq 1$, then $\mathbf{w}_2 = \frac{1}{\lambda-1} \mathbf{G}^\top \mathbf{w}_1$. Substituting in equation (2),

$$dL(\mathbf{w}_1, d\mathbf{w}_1) = -2(\mathbf{w}_1^\top \mathbf{H} + \mathbf{w}_1^\top \mathbf{G} \mathbf{G}^\top + \frac{1}{\lambda-1} \mathbf{w}_1^\top \mathbf{G} \mathbf{G}^\top - \lambda \mathbf{w}_1^\top) d\mathbf{w}_1.$$

The first order holds if and only if, $(\mathbf{H} + \frac{\lambda}{\lambda-1} \mathbf{G} \mathbf{G}^\top) \mathbf{w}_1 = \lambda \mathbf{w}_1$. \square

1080 D.3 PROPOSITION 1
1081

1082 **Proposition 1** (Behaviour after measurement). *If $(\mathbf{l}, \mathbf{x}, \mathbf{y})$ are distributed according to figure 2. Then*
1083 *the following properties hold,*

1084 CALIBRATED INDEPENDENCE: $S(\mathbf{x} | \mathbf{l}) = 0 \iff \mathbf{x} \perp \mathbf{l}$ and $S(\mathbf{y}; \mathbf{l} | \mathbf{x}) = 0 \iff \mathbf{y} | \mathbf{x} \perp \mathbf{l} | \mathbf{x}$.

1085 CHAIN RULE OF INFORMATION: $S(\mathbf{x}; \mathbf{y} | \mathbf{l}) = S(\mathbf{x}; \mathbf{l}) + S(\mathbf{y}; \mathbf{l} | \mathbf{x})$.

1086 STATE INDEPENDENT BAYES: $S(\mathbf{x}; \mathbf{l}) = S(\mathbf{y}; \mathbf{l} | \mathbf{x}) = 0 \implies S(\mathbf{y}; \mathbf{l}) = S(\mathbf{x}; \mathbf{l} | \mathbf{y}) = 0$.

1087

1088 *Proof.* (i) Follows from Jensen's inequality. (ii) from theorem 4. (iii) from proposition 10. \square
1089

1090 D.3.1 SYMMETRIC DEPENDENCE ON THE STATE
1091

1092 **Definition and construction**

1093

1094 UNCONDITIONAL DEPENDENCE TO STATES Two random variables are independent if and only if their
1095 joint distribution factors into the product of their marginals. This, in conjunction with definition 9
1096 motivates the following definition,

1097 **Definition 5** (Unconditional f -dependence to state). The *unconditional f -dependence to state* is
1098 defined as,

$$1099 I_f(\mathbf{x}; \mathbf{l}) := D_f(\mathbb{P}(\mathbf{x}, \mathbf{l}) || \mathbb{P}(\mathbf{x})\mathbb{P}(\mathbf{l})).$$

1100

1101

1102 The next proposition expresses the unconditional f -dependence to states,

1103 **Proposition 7.** *The unconditional f -dependence to state is equal to,*

1104

$$1105 I_f(\mathbf{x}; \mathbf{l}) = \mathbb{E} \left[f \left(\frac{\mathbb{P}(\mathbf{x} | \mathbf{l})}{\mathbb{P}(\mathbf{x})} \right) \mathbb{P}(\mathbf{x} | \mathbf{l}') \right].$$

1106

1107 *Where \mathbf{l}' is an independent copy of \mathbf{l} .*

1108

1109 *Proof.* From definition 5,

1110

$$1111 I_f = D_f(\mathbb{P}(\mathbf{x}, \mathbf{l}) || \mathbb{P}(\mathbf{x})\mathbb{P}(\mathbf{l})),$$

$$1112 = \sum_{i=1}^{|\mathcal{I}|} \left(\int f \left(\frac{\mathbb{P}(\mathbf{x}, \mathbf{l} = i)}{\mathbb{P}(\mathbf{x})\mathbb{P}(\mathbf{l} = i)} \right) \mathbb{P}(\mathbf{x}) \, d\mathbf{x} \right) \mathbb{P}(\mathbf{l} = i),$$

$$1113 = \sum_{i=1}^{|\mathcal{I}|} \left(\int f \left(\frac{\mathbb{P}(\mathbf{x}, \mathbf{l} = i)}{\mathbb{P}(\mathbf{x})\mathbb{P}(\mathbf{l} = i)} \right) \sum_{j=1}^{|\mathcal{I}|} \mathbb{P}(\mathbf{x} | \mathbf{l} = j) \, d\mathbf{x} \right) \mathbb{P}(\mathbf{l} = i),$$

$$1114 = \sum_{1 \leq i, j \leq |\mathcal{I}|} \mathbb{P}(\mathbf{l} = j)\mathbb{P}(\mathbf{l} = i) \left(\int f \left(\frac{\mathbb{P}(\mathbf{x}, \mathbf{l} = i)}{\mathbb{P}(\mathbf{x})\mathbb{P}(\mathbf{l} = i)} \right) \mathbb{P}(\mathbf{x} | \mathbf{l} = j) \, d\mathbf{x} \right),$$

$$1115 = \sum_{1 \leq i, j \leq |\mathcal{I}|} \mathbb{P}(\mathbf{l} = j)\mathbb{P}(\mathbf{l} = i) \left(\int f \left(\frac{\mathbb{P}(\mathbf{x} | \mathbf{l} = i)}{\mathbb{P}(\mathbf{x})} \right) \mathbb{P}(\mathbf{x} | \mathbf{l} = j) \, d\mathbf{x} \right).$$

1123

1124 \square

1125

1126 The unconditional f -dependence is determined by convex combinations of terms of the form,
1127 $G_f(i, j) := \int f \left(\frac{\mathbb{P}(\mathbf{x} | \mathbf{l} = i)}{\mathbb{P}(\mathbf{x})} \right) \mathbb{P}(\mathbf{x} | \mathbf{l} = j) \, d\mathbf{x}$. If we consider the f -dependence with generator
1128 $t \mapsto f(t) = t \ln t$. The unconditional dependence to state is exactly the mutual information of \mathbf{x}
1129 and \mathbf{l} .

1130 **Corollary 5.** *Let $f(t) = t \ln t$, the Kullback-Leibler divergence generator. Then,*

1131

$$1132 I_{KL}(\mathbf{x}; \mathbf{l}) = I(\mathbf{x}; \mathbf{l}) = H[\mathbb{P}(\mathbf{x})] - \mathbb{E}[H[\mathbb{P}(\mathbf{x} | \mathbf{l})]].$$

1133

The mutual information of \mathbf{x} and \mathbf{l} .

1134
1135
1136
1137
1138
1139
1140
1141
1142
1143
1144
1145
1146
1147
1148
1149
1150
1151
1152
1153
1154
1155
1156
1157
1158
1159
1160
1161
1162
1163
1164
1165
1166
1167
1168
1169
1170
1171
1172
1173
1174
1175
1176
1177
1178
1179
1180
1181
1182
1183
1184
1185
1186
1187

Proof. Indeed in this setting,

$$\begin{aligned}
G_f(i, j) &= \left(\frac{\mathbb{P}(\mathbf{x} | \mathbf{I} = i)}{\mathbb{P}(\mathbf{x})} \right) \ln \left(\frac{\mathbb{P}(\mathbf{x} | \mathbf{I} = i)}{\mathbb{P}(\mathbf{x})} \right) \mathbb{P}(\mathbf{x} | \mathbf{I} = j), \\
&= \left(\frac{\mathbb{P}(\mathbf{x} | \mathbf{I} = i)}{\mathbb{P}(\mathbf{x})} \right) (\ln \mathbb{P}(\mathbf{x} | \mathbf{I} = i) - \ln \mathbb{P}(\mathbf{x})) \mathbb{P}(\mathbf{x} | \mathbf{I} = j), \\
&= \left(\frac{\mathbb{P}(\mathbf{x} | \mathbf{I} = i) \mathbb{P}(\mathbf{x} | \mathbf{I} = j)}{\mathbb{P}(\mathbf{x})} \right) (\ln \mathbb{P}(\mathbf{x} | \mathbf{I} = i) - \ln \mathbb{P}(\mathbf{x})), \\
&= \left(\ln \mathbb{P}(\mathbf{x} | \mathbf{I} = i) \left(\frac{\mathbb{P}(\mathbf{x} | \mathbf{I} = i) \mathbb{P}(\mathbf{x} | \mathbf{I} = j)}{\mathbb{P}(\mathbf{x})} \right) - \ln \mathbb{P}(\mathbf{x}) \left(\frac{\mathbb{P}(\mathbf{x} | \mathbf{I} = i) \mathbb{P}(\mathbf{x} | \mathbf{I} = j)}{\mathbb{P}(\mathbf{x})} \right) \right).
\end{aligned}$$

Summing over j with weights $\mathbb{P}(\mathbf{I} = j)$,

$$\begin{aligned}
\sum_{j=1}^{|\mathcal{I}|} \mathbb{P}(\mathbf{I} = j) G_f(i, j) &= \ln \mathbb{P}(\mathbf{x} | \mathbf{I} = i) \left(\frac{\mathbb{P}(\mathbf{x} | \mathbf{I} = i) \sum_{j=1}^{|\mathcal{I}|} \mathbb{P}(\mathbf{I} = j) \mathbb{P}(\mathbf{x} | \mathbf{I} = j)}{\mathbb{P}(\mathbf{x})} \right) \\
&\quad - \ln \mathbb{P}(\mathbf{x}) \left(\frac{\mathbb{P}(\mathbf{x} | \mathbf{I} = i) \sum_{j=1}^{|\mathcal{I}|} \mathbb{P}(\mathbf{I} = j) \mathbb{P}(\mathbf{x} | \mathbf{I} = j)}{\mathbb{P}(\mathbf{x})} \right), \\
&= \ln \mathbb{P}(\mathbf{x} | \mathbf{I} = i) \left(\frac{\mathbb{P}(\mathbf{x} | \mathbf{I} = i) \mathbb{P}(\mathbf{x})}{\mathbb{P}(\mathbf{x})} \right) - \ln \mathbb{P}(\mathbf{x}) \left(\frac{\mathbb{P}(\mathbf{x} | \mathbf{I} = i) \mathbb{P}(\mathbf{x})}{\mathbb{P}(\mathbf{x})} \right), \\
&= \ln \mathbb{P}(\mathbf{x} | \mathbf{I} = i) \mathbb{P}(\mathbf{x} | \mathbf{I} = i) - \ln \mathbb{P}(\mathbf{x}) \mathbb{P}(\mathbf{x} | \mathbf{I} = i).
\end{aligned}$$

Summing over i with weights $\mathbb{P}(\mathbf{I} = i)$,

$$\sum_{1 \leq i, j \leq |\mathcal{I}|} \mathbb{P}(\mathbf{I} = j) G_f(i, j) = \sum_{i=1}^{|\mathcal{I}|} \mathbb{P}(\mathbf{I} = i) \ln \mathbb{P}(\mathbf{x} | \mathbf{I} = i) \mathbb{P}(\mathbf{x} | \mathbf{I} = i) - \ln \mathbb{P}(\mathbf{x}) \mathbb{P}(\mathbf{x}).$$

Integrating with respect to $\mathbb{P}(\mathbf{x})$,

$$\begin{aligned}
U_f &= \sum_{i=1}^{|\mathcal{I}|} \mathbb{P}(\mathbf{I} = i) \int \ln \mathbb{P}(\mathbf{x} | \mathbf{I} = i) \mathbb{P}(\mathbf{x} | \mathbf{I} = i) \, d\mathbf{x} - \int \ln \mathbb{P}(\mathbf{x}) \mathbb{P}(\mathbf{x}) \, d\mathbf{x}, \\
&= H[\mathbb{P}(\mathbf{x})] - \sum_{i=1}^{|\mathcal{I}|} \mathbb{P}(\mathbf{I} = i) H[\mathbb{P}(\mathbf{x} | \mathbf{I} = i)] = H[\mathbb{P}(\mathbf{x})] - \mathbb{E}[H[\mathbb{P}(\mathbf{x} | \mathbf{I})]] = I(\mathbf{x}; \mathbf{I}).
\end{aligned}$$

□

Now we consider the reverse Kullback-Leibler divergence.

Corollary 6. Let $f(t) = -\ln t$, reverse Kullback-Leibler divergence generator. Then,

$$I_{LK}(\mathbf{x}; \mathbf{I}) = \mathbb{E}[H[\mathbb{P}(\mathbf{x} | \mathbf{I}), \mathbb{P}(\mathbf{x} | \mathbf{I}')] - H[\mathbb{P}(\mathbf{x})]].$$

Where \mathbf{I}' is independent copy of \mathbf{I} .

Proof. Indeed in this setting, $-G_f(i, j) = \ln \left(\frac{\mathbb{P}(\mathbf{x} | \mathbf{I} = i)}{\mathbb{P}(\mathbf{x})} \right) \mathbb{P}(\mathbf{x} | \mathbf{I} = j) = (\ln \mathbb{P}(\mathbf{x} | \mathbf{I} = i) - \ln \mathbb{P}(\mathbf{x})) \mathbb{P}(\mathbf{x} | \mathbf{I} = j)$. Summing over j with weights $\mathbb{P}(\mathbf{I} = j)$,

$$-\sum_{j=1}^{|\mathcal{I}|} \mathbb{P}(\mathbf{I} = j) G_f(i, j) = \sum_{j=1}^{|\mathcal{I}|} \mathbb{P}(\mathbf{I} = j) \ln \mathbb{P}(\mathbf{x} | \mathbf{I} = i) \mathbb{P}(\mathbf{x} | \mathbf{I} = j) - \ln \mathbb{P}(\mathbf{x}) \mathbb{P}(\mathbf{x}).$$

Integrating with respect to \mathbf{x} , and summing over i with weights $\mathbb{P}(\mathbf{I} = i)$,

$$\begin{aligned}
-U_f(\mathbf{x}; \mathbf{I}) &= \sum_{1 \leq i, j \leq |\mathcal{I}|} \mathbb{P}(\mathbf{I} = i) \mathbb{P}(\mathbf{I} = j) \int \ln \mathbb{P}(\mathbf{x} | \mathbf{I} = i) \mathbb{P}(\mathbf{x} | \mathbf{I} = j) \, d\mathbf{x} - \int \ln \mathbb{P}(\mathbf{x}) \mathbb{P}(\mathbf{x}) \, d\mathbf{x}, \\
&= H[\mathbb{P}(\mathbf{x})] - \sum_{1 \leq i, j \leq |\mathcal{I}|} \mathbb{P}(\mathbf{I} = i) \mathbb{P}(\mathbf{I} = j) H[\mathbb{P}(\mathbf{x} | \mathbf{I} = i), \mathbb{P}(\mathbf{x} | \mathbf{I} = j)], \\
&= H[\mathbb{P}(\mathbf{x})] - \mathbb{E}[H[\mathbb{P}(\mathbf{x} | \mathbf{I}), \mathbb{P}(\mathbf{x} | \mathbf{I}')]].
\end{aligned}$$

And hence, $U_f(\mathbf{x}; \mathbf{I}) = \mathbb{E}[H[\mathbb{P}(\mathbf{x} | \mathbf{I}), \mathbb{P}(\mathbf{x} | \mathbf{I}')] - H[\mathbb{P}(\mathbf{x})]]$.

□

1188 By exploiting exploiting the properties of f -divergences we can readily express the analogue for the
 1189 symmetrized Kullblack-Leilber divergence,

1190 **Corollary 7.** When $f(t) = (t - 1) \ln t$, the Jeffrey's divergence generator,

$$1192 I_J(\mathbf{x}; \mathbf{l}) = \mathbb{E} \left[KL \left(\mathbb{P}(\mathbf{x} | \mathbf{l}) \parallel \mathbb{P}(\mathbf{x} | \mathbf{l}') \right) \right].$$

1194 Where \mathbf{l}' is an independent copy of \mathbf{l} .

1196 *Proof.* Let $g(t) = t \ln t$ and $h(t) = -\ln t$ denote the respective generators of the Kullblack-Leibler
 1197 and reverse Kullblack-Leibler divergences. We can write, $f(t) = g(t) + h(t)$. Hence by Proposition 5,
 1198 $U_f = S_g + S_h$. By 5 and 6,

$$1200 I_J(\mathbf{x}; \mathbf{l}) = \mathbb{E}[H[\mathbb{P}(\mathbf{x} | \mathbf{l}), \mathbb{P}(\mathbf{x} | \mathbf{l}')] - H[\mathbb{P}(\mathbf{x})] + H[\mathbb{P}(\mathbf{x})] - \mathbb{E}[H[\mathbb{P}(\mathbf{x} | \mathbf{l})]],$$

$$1202 = \mathbb{E}[H[\mathbb{P}(\mathbf{x} | \mathbf{l}), \mathbb{P}(\mathbf{x} | \mathbf{l}')] - \mathbb{E}[H[\mathbb{P}(\mathbf{x} | \mathbf{l})]] = \mathbb{E} \left[KL \left(\mathbb{P}(\mathbf{x} | \mathbf{l}) \parallel \mathbb{P}(\mathbf{x} | \mathbf{l}') \right) \right].$$

1204 □

1205 Which we can express in terms of symmetrized Kullblack-Leibler divergences pairs sampled without
 1206 replacement

1208 **Corollary 8.**

$$1209 I_J(\mathbf{x}; \mathbf{l}) = \frac{1}{2} \mathbb{E} \left[D_J \left(\mathbb{P}(\mathbf{x} | \mathbf{l}) \parallel \mathbb{P}(\mathbf{x} | \mathbf{l}') \right) \right].$$

1211 Where \mathbf{l} is an independent copy of \mathbf{l}' .

1213 *Proof.* Follows from the fact that Jeffreys' divergence and symmetric and corollary 7. □

1215 **CONDITIONAL** We can now move to characterizing the conditional dependence. \mathbf{y} and \mathbf{l} are in-
 1216 dependent given \mathbf{x} if and and only, $\mathbb{P}(\mathbf{y}, \mathbf{l} | \mathbf{x}) = \mathbb{P}(\mathbf{y} | \mathbf{x})\mathbb{P}(\mathbf{l} | \mathbf{x})$. This prompts the following
 1217 definition.

1218 **Definition 6** (conditional f -dependence to state). $I_f(\mathbf{y}; \mathbf{l} | \mathbf{x}) :=$
 1219 $D_f(\mathbb{P}(\mathbf{y}, \mathbf{l} | \mathbf{x}) \parallel \mathbb{P}(\mathbf{y} | \mathbf{x})\mathbb{P}(\mathbf{l} | \mathbf{x}) | \mathbb{P}(\mathbf{x}))$.

1220 The following proposition yields a handful of useful expression for the conditional dependence.

1221 **Corollary 9.** Let $\mathbf{x} \mapsto U_f(\mathbf{x}) = D_f(\mathbb{P}(\mathbf{y}, \mathbf{l} | \mathbf{x}) \parallel \mathbb{P}(\mathbf{y} | \mathbf{x})\mathbb{P}(\mathbf{l} | \mathbf{x}))$. The following are all equal to
 1222 $I_f(\mathbf{y}; \mathbf{l} | \mathbf{x})$.

- 1224 (i) $\int U_f(\mathbf{x})\mathbb{P}(\mathbf{x}) \, d\mathbf{x}$,
- 1225 (ii) $\int \left(\sum_{i=1}^{|\mathcal{I}|} \int f \left(\frac{\mathbb{P}(\mathbf{y}, \mathbf{l}=i | \mathbf{x})}{\mathbb{P}(\mathbf{y} | \mathbf{x})\mathbb{P}(\mathbf{l}=i | \mathbf{x})} \right) \mathbb{P}(\mathbf{y} | \mathbf{x})\mathbb{P}(\mathbf{l} = i | \mathbf{x}) \, d\mathbf{y} \right) \mathbb{P}(\mathbf{x}) \, d\mathbf{x}$,
- 1226 (iii) $\int \sum_{1 \leq i, j \leq |\mathcal{I}|} \mathbb{P}(\mathbf{x})\mathbb{P}(\mathbf{l} = i | \mathbf{x})\mathbb{P}(\mathbf{l} = j | \mathbf{x}) \left(\int f \left(\frac{\mathbb{P}(\mathbf{y} | \mathbf{l}=i, \mathbf{x})}{\mathbb{P}(\mathbf{y} | \mathbf{x})} \right) \mathbb{P}(\mathbf{y} | \mathbf{l} = j, \mathbf{x}) \, d\mathbf{y} \right) \, d\mathbf{x}$,
- 1227 (iv) $\sum_{1 \leq i, j, k \leq |\mathcal{I}|} \mathbb{P}(\mathbf{l} = k) \int \mathbb{P}(\mathbf{l} = i | \mathbf{x})\mathbb{P}(\mathbf{l} = j | \mathbf{x}) \left(\int f \left(\frac{\mathbb{P}(\mathbf{y} | \mathbf{l}=i, \mathbf{x})}{\mathbb{P}(\mathbf{y} | \mathbf{x})} \right) \mathbb{P}(\mathbf{y} | \mathbf{l} = j, \mathbf{x}) \, d\mathbf{y} \right) \mathbb{P}(\mathbf{x} |$
 1228 $\mathbf{l} = k) \, d\mathbf{x}$.

1232 *Proof.* This in turn can be expressed as,

$$1233 I_f(\mathbf{y}; \mathbf{l} | \mathbf{x}) = \int \sum_{1 \leq i, j \leq |\mathcal{I}|} \mathbb{P}(\mathbf{l} = i | \mathbf{x})\mathbb{P}(\mathbf{l} = j | \mathbf{x}) \left(\int f \left(\frac{\mathbb{P}(\mathbf{y} | \mathbf{l} = i, \mathbf{x})}{\mathbb{P}(\mathbf{y} | \mathbf{x})} \right) \mathbb{P}(\mathbf{y} | \mathbf{l} = j, \mathbf{x}) \, d\mathbf{y} \right) \mathbb{P}(\mathbf{x}) \, d\mathbf{x},$$

$$1234 = \int \sum_{1 \leq i, j \leq |\mathcal{I}|} \mathbb{P}(\mathbf{x})\mathbb{P}(\mathbf{l} = i | \mathbf{x})\mathbb{P}(\mathbf{l} = j | \mathbf{x}) \left(\int f \left(\frac{\mathbb{P}(\mathbf{y} | \mathbf{l} = i, \mathbf{x})}{\mathbb{P}(\mathbf{y} | \mathbf{x})} \right) \mathbb{P}(\mathbf{y} | \mathbf{l} = j, \mathbf{x}) \, d\mathbf{y} \right) \, d\mathbf{x}, \quad (iii).$$

1239 By Bayes theorem we have,

$$1241 \mathbb{P}(\mathbf{l} = i | \mathbf{x}) = \frac{\mathbb{P}(\mathbf{l} = i)\mathbb{P}(\mathbf{x} | \mathbf{l} = i)}{\mathbb{P}(\mathbf{x})}.$$

1242
1243
1244
1245
1246
1247
1248
1249
1250
1251
1252
1253
1254
1255
1256
1257
1258
1259
1260
1261
1262
1263
1264
1265
1266
1267
1268
1269
1270
1271
1272
1273
1274
1275
1276
1277
1278
1279
1280
1281
1282
1283
1284
1285
1286
1287
1288
1289
1290
1291
1292
1293
1294
1295

Hence,

$$I_f(\mathbf{y}; \mathbf{l} | \mathbf{x}) = \sum_{1 \leq i, j \leq |\mathcal{L}|} \mathbb{P}(\mathbf{l} = i) \int \mathbb{P}(\mathbf{x} | \mathbf{l} = i) \mathbb{P}(\mathbf{l} = j | \mathbf{x}) \left(\int f \left(\frac{\mathbb{P}(\mathbf{y} | \mathbf{l} = i, \mathbf{x})}{\mathbb{P}(\mathbf{y} | \mathbf{x})} \right) \mathbb{P}(\mathbf{y} | \mathbf{l} = j, \mathbf{x}) \, d\mathbf{y} \right) \, d\mathbf{x}.$$

Moreover,

$$\begin{aligned} I_f(\mathbf{y}; \mathbf{l} | \mathbf{x}) &= \int \sum_{1 \leq i, j \leq |\mathcal{L}|} \mathbb{P}(\mathbf{l} = i | \mathbf{x}) \mathbb{P}(\mathbf{l} = j | \mathbf{x}) \left(\int f \left(\frac{\mathbb{P}(\mathbf{y} | \mathbf{l} = i, \mathbf{x})}{\mathbb{P}(\mathbf{y} | \mathbf{x})} \right) \mathbb{P}(\mathbf{y} | \mathbf{l} = j, \mathbf{x}) \, d\mathbf{y} \right) \mathbb{P}(\mathbf{x}) \, d\mathbf{x}, \\ &= \sum_{1 \leq i, j, k \leq |\mathcal{L}|} \mathbb{P}(\mathbf{l} = k) \int \mathbb{P}(\mathbf{l} = i | \mathbf{x}) \mathbb{P}(\mathbf{l} = j | \mathbf{x}) \left(\int f \left(\frac{\mathbb{P}(\mathbf{y} | \mathbf{l} = i, \mathbf{x})}{\mathbb{P}(\mathbf{y} | \mathbf{x})} \right) \mathbb{P}(\mathbf{y} | \mathbf{l} = j, \mathbf{x}) \, d\mathbf{y} \right) \mathbb{P}(\mathbf{x} | \mathbf{l} = k) \, d\mathbf{x}. \quad (iv) \end{aligned}$$

□

We can now readily express the LK , KL and J conditional sensitivities,

Corollary 10. $I_{KL}(\mathbf{y}; \mathbf{l} | \mathbf{x}) = H(\mathbf{y} | \mathbf{x}) - \mathbb{E} [H[\mathbb{P}(\mathbf{y} | \mathbf{l}, \mathbf{x})]]$.

Proof. By corollary 5, $U_{KL}(\mathbf{x}) = H[\mathbb{P}(\mathbf{y} | \mathbf{x})] - \mathbb{E} [H[\mathbb{P}(\mathbf{y} | \mathbf{l}, \mathbf{x}) | \mathbf{x}]]$. Hence,

$$\begin{aligned} I_{KL}(\mathbf{y}; \mathbf{l} | \mathbf{x}) &= \mathbb{E} [H[\mathbb{P}(\mathbf{y} | \mathbf{x})]] - \mathbb{E} [\mathbb{E} [H[\mathbb{P}(\mathbf{y} | \mathbf{l}, \mathbf{x}) | \mathbf{x}]]], \\ &= H(\mathbf{y} | \mathbf{x}) - \mathbb{E} [H[\mathbb{P}(\mathbf{y} | \mathbf{l}, \mathbf{x})]]. \end{aligned}$$

□

Proposition 8. $I_{LK}(\mathbf{y}; \mathbf{l} | \mathbf{x}) = \mathbb{E} [H[\mathbb{P}(\mathbf{y} | \mathbf{l}, \mathbf{x}), \mathbb{P}(\mathbf{y} | \mathbf{l}', \mathbf{x}) | \mathbf{x}]] - \mathbb{E} [H[\mathbb{P}(\mathbf{y} | \mathbf{x})]]$.

Proof. By corollary 5 $U_{LK}(\mathbf{x}) = \mathbb{E} [H[\mathbb{P}(\mathbf{y} | \mathbf{l}, \mathbf{x}), \mathbb{P}(\mathbf{y} | \mathbf{l}', \mathbf{x}) | \mathbf{x}]] - H[\mathbb{P}(\mathbf{y} | \mathbf{x})]$. Hence,

$$\begin{aligned} I_{LK}(\mathbf{y}; \mathbf{l} | \mathbf{x}) &= \mathbb{E} [\mathbb{E} [H[\mathbb{P}(\mathbf{y} | \mathbf{l}, \mathbf{x}), \mathbb{P}(\mathbf{y} | \mathbf{l}', \mathbf{x}) | \mathbf{x}]]] - \mathbb{E} [H[\mathbb{P}(\mathbf{y} | \mathbf{x})]], \\ &= \mathbb{E} [H[\mathbb{P}(\mathbf{y} | \mathbf{l}, \mathbf{x}), \mathbb{P}(\mathbf{y} | \mathbf{l}', \mathbf{x}) | \mathbf{x}]] - \mathbb{E} [H[\mathbb{P}(\mathbf{y} | \mathbf{x})]]. \end{aligned}$$

□

Corollary 11. $I_J(\mathbf{y}; \mathbf{l} | \mathbf{x}) = \mathbb{E} \left[KL \left(\mathbb{P}(\mathbf{y} | \mathbf{l}, \mathbf{x}) \parallel \mathbb{P}(\mathbf{y} | \mathbf{l}', \mathbf{x}) \mid \mathbb{P}(\mathbf{x}) \right) \right]$, Where \mathbf{l}' is independent copy of \mathbf{l} .

As in corollary 7 we can express the conditional dependence to state with respect to the symmetrized Kullback-Leibler divergence.

Corollary 12. $I_J(\mathbf{y}; \mathbf{l} | \mathbf{x}) = \frac{1}{2} \mathbb{E} \left[D_J \left(\mathbb{P}(\mathbf{y} | \mathbf{x}, \mathbf{l}) \parallel \mathbb{P}(\mathbf{y} | \mathbf{x}, \mathbf{l}') \mid \mathbb{P}(\mathbf{x}) \right) \right]$, Where \mathbf{l}' is an independent copy of \mathbf{l} .

Chain rule of information We can now establish the information Bayes theorem for the I_J . First we need to establish the chain rule for the Kullback-Leibler divergence

Lemma 3 (Chain rule for the Kullback-Leibler divergence). $KL(\mathbb{P}_{\mathbf{xy}} \parallel \mathbb{Q}_{\mathbf{xy}}) = KL(\mathbb{P}_{\mathbf{y|x}} \parallel \mathbb{Q}_{\mathbf{y|x}} \mid \mathbb{P}_{\mathbf{x}}) + KL(\mathbb{P}_{\mathbf{x}} \parallel \mathbb{Q}_{\mathbf{x}})$.

Proof. (Polyanskiy & Wu, 2024)

□

Proposition 9.

Proof. Start by decomposing the likelihood ratio,

$$\frac{\mathbb{P}_{\mathbf{xy}|\mathbf{l}}}{\mathbb{P}_{\mathbf{xy}} \mathbb{P}_{\mathbf{l}}} = \frac{\mathbb{P}_{\mathbf{y}|\mathbf{l}|\mathbf{x}}}{\mathbb{P}_{\mathbf{y}|\mathbf{x}} \mathbb{P}_{\mathbf{l}|\mathbf{x}}} \frac{\mathbb{P}_{\mathbf{x}|\mathbf{l}}}{\mathbb{P}_{\mathbf{x}} \mathbb{P}_{\mathbf{l}}}. \quad (3)$$

1296 Applying the function $t \mapsto t \log t$ on each side of equation (3),
 1297

$$1298 \frac{\mathbb{P}_{\mathbf{xy}|\mathbf{l}}}{\mathbb{P}_{\mathbf{xy}}\mathbb{P}_{\mathbf{l}}} \log \left(\frac{\mathbb{P}_{\mathbf{xy}|\mathbf{l}}}{\mathbb{P}_{\mathbf{xy}}\mathbb{P}_{\mathbf{l}}} \right) = \frac{\mathbb{P}_{\mathbf{y}|\mathbf{x}}}{\mathbb{P}_{\mathbf{y}|\mathbf{x}}\mathbb{P}_{\mathbf{l}|\mathbf{x}}} \frac{\mathbb{P}_{\mathbf{x}|\mathbf{l}}}{\mathbb{P}_{\mathbf{x}}\mathbb{P}_{\mathbf{l}}} \log \left(\frac{\mathbb{P}_{\mathbf{y}|\mathbf{x}}}{\mathbb{P}_{\mathbf{y}|\mathbf{x}}\mathbb{P}_{\mathbf{l}|\mathbf{x}}} \frac{\mathbb{P}_{\mathbf{x}|\mathbf{l}}}{\mathbb{P}_{\mathbf{x}}\mathbb{P}_{\mathbf{l}}} \right),$$

$$1300 = \frac{\mathbb{P}_{\mathbf{y}|\mathbf{x}}}{\mathbb{P}_{\mathbf{y}|\mathbf{x}}\mathbb{P}_{\mathbf{l}|\mathbf{x}}} \frac{\mathbb{P}_{\mathbf{x}|\mathbf{l}}}{\mathbb{P}_{\mathbf{x}}\mathbb{P}_{\mathbf{l}}} \left(\log \left(\frac{\mathbb{P}_{\mathbf{y}|\mathbf{x}}}{\mathbb{P}_{\mathbf{y}|\mathbf{x}}\mathbb{P}_{\mathbf{l}|\mathbf{x}}} \right) + \log \left(\frac{\mathbb{P}_{\mathbf{x}|\mathbf{l}}}{\mathbb{P}_{\mathbf{x}}\mathbb{P}_{\mathbf{l}}} \right) \right).$$

1303 Taking the Expectation with respect to $\mathbb{P}_{\mathbf{xy}}\mathbb{P}_{\mathbf{l}}$ the left hand side is equal to,
 1304

$$1305 \mathbb{E}_{\mathbb{P}_{\mathbf{xy}}\mathbb{P}_{\mathbf{l}}} \left[\frac{\mathbb{P}_{\mathbf{xy}|\mathbf{l}}}{\mathbb{P}_{\mathbf{xy}}\mathbb{P}_{\mathbf{l}}} \log \left(\frac{\mathbb{P}_{\mathbf{xy}|\mathbf{l}}}{\mathbb{P}_{\mathbf{xy}}\mathbb{P}_{\mathbf{l}}} \right) \right] = \mathbb{E}_{\mathbb{P}_{\mathbf{xy}|\mathbf{l}}} \left[\log \left(\frac{\mathbb{P}_{\mathbf{xy}|\mathbf{l}}}{\mathbb{P}_{\mathbf{xy}}\mathbb{P}_{\mathbf{l}}} \right) \right] = I_{KL}([\mathbf{x} \mathbf{y}]; \mathbf{l}).$$

1307 The integration measure on the right hand side simplifies to,
 1308

$$1309 \mathbb{P}_{\mathbf{xy}}\mathbb{P}_{\mathbf{l}} \frac{\mathbb{P}_{\mathbf{y}|\mathbf{x}}}{\mathbb{P}_{\mathbf{y}|\mathbf{x}}\mathbb{P}_{\mathbf{l}|\mathbf{x}}} \frac{\mathbb{P}_{\mathbf{x}|\mathbf{l}}}{\mathbb{P}_{\mathbf{x}}\mathbb{P}_{\mathbf{l}}} = \frac{\mathbb{P}_{\mathbf{y}|\mathbf{x}}}{\mathbb{P}_{\mathbf{l}|\mathbf{x}}} \mathbb{P}_{\mathbf{x}|\mathbf{l}} = \frac{\mathbb{P}_{\mathbf{y}|\mathbf{x}}}{\mathbb{P}_{\mathbf{l}|\mathbf{x}}} \mathbb{P}_{\mathbf{l}|\mathbf{x}} \mathbb{P}_{\mathbf{x}} = \mathbb{P}_{\mathbf{xy}|\mathbf{l}}.$$

1312 Distributing over the logarithms, the right hand side is equal, $I_{KL}(\mathbf{y}; \mathbf{l} | \mathbf{x}) + I(\mathbf{x}; \mathbf{l})$. Similarly by
 1313 flipping \mathbf{x} and \mathbf{y} , $I_{KL}([\mathbf{x} \mathbf{y}]; \mathbf{l}) = I_{KL}(\mathbf{x}; \mathbf{l} | \mathbf{y}) + I(\mathbf{y}; \mathbf{l})$. Hence, $I_{KL}(\mathbf{y}; \mathbf{l} | \mathbf{x}) + I_{KL}(\mathbf{x}; \mathbf{l} | \mathbf{y}) + I_{KL}(\mathbf{x}; \mathbf{l} | \mathbf{y}) + I_{KL}(\mathbf{y}; \mathbf{l})$. \square
 1314

1315 We can now prove the assertion that it is enough for use to consider unconditional and conditional
 1316 sensitivities to states.
 1317

1318 **Theorem 4.**

$$1319 S(\mathbf{x}; \mathbf{l} | \mathbf{y}) + S(\mathbf{y} | \mathbf{l}) = S([\mathbf{x} \mathbf{y}] | \mathbf{l}) = S([\mathbf{y} \mathbf{x}] | \mathbf{l}) = S(\mathbf{y}; \mathbf{l} | \mathbf{x}) + S(\mathbf{x} | \mathbf{l})$$

1322 *Proof.* By corollary 7, $S([\mathbf{x} \mathbf{y}] | \mathbf{l}) = \mathbb{E}[KL(\mathbb{P}_{\mathbf{xy}|\mathbf{l}} || \mathbb{P}_{\mathbf{xy}|\mathbf{l}'})]$. Where \mathbf{l}' is an independent copy of
 1323 \mathbf{l} . By lemma 3, $KL(\mathbb{P}_{\mathbf{xy}|\mathbf{l}} || \mathbb{P}_{\mathbf{xy}|\mathbf{l}'}) = KL(\mathbb{P}_{\mathbf{y}|\mathbf{x},\mathbf{l}} || \mathbb{P}_{\mathbf{y}|\mathbf{x},\mathbf{l}'}) + KL(\mathbb{P}_{\mathbf{x}|\mathbf{l}} || \mathbb{P}_{\mathbf{x}|\mathbf{l}'})$. Taking the
 1324 expectation with respect to $(\mathbf{l}, \mathbf{l}')$, $S([\mathbf{x} \mathbf{y}] | \mathbf{l}) = \mathbb{E}[KL(\mathbb{P}_{\mathbf{y}|\mathbf{x},\mathbf{l}} || \mathbb{P}_{\mathbf{y}|\mathbf{x},\mathbf{l}'})] + \mathbb{E}[KL(\mathbb{P}_{\mathbf{x}|\mathbf{l}} || \mathbb{P}_{\mathbf{x}|\mathbf{l}'})]$.
 1325 Hence, $S([\mathbf{x} \mathbf{y}] | \mathbf{l}) = S(\mathbf{y}; \mathbf{l} | \mathbf{x}) + S(\mathbf{x}; \mathbf{l} | \mathbf{l})$. Switching the roles of \mathbf{x} and \mathbf{y} , by symmetry we establish
 1326 the information Bayes theorem, $S(\mathbf{x}; \mathbf{l} | \mathbf{y}) + S(\mathbf{y}; \mathbf{l}) = S([\mathbf{x} \mathbf{y}] | \mathbf{l}) = S([\mathbf{y} \mathbf{x}] | \mathbf{l}) = S(\mathbf{y}; \mathbf{l} | \mathbf{x}) + S(\mathbf{x}; \mathbf{l} | \mathbf{l})$. \square
 1327
 1328
 1329

1330 **State Independent Bayes**

1331 **Proposition 10** (State Independent Bayes). (i) $S(\mathbf{y}; \mathbf{l}) = 0 \wedge S(\mathbf{x}; \mathbf{l} | \mathbf{y}) = 0 \implies S(\mathbf{y}; \mathbf{l} | \mathbf{x}) = 0$ (Bayes).
 1332

1333 (ii) $S(\mathbf{y}; \mathbf{l}) = 0 \wedge S(\mathbf{x}; \mathbf{l} | \mathbf{y}) = 0 \implies S(\mathbf{x}; \mathbf{l}) = 0$ (Marginalization)
 1334

1335 *Proof.* $S(\mathbf{y}; \mathbf{l} | \mathbf{x}) = 0 \wedge S(\mathbf{x}; \mathbf{l}) = 0 \implies S(\mathbf{y}; \mathbf{l}) = 0 \wedge S(\mathbf{x}; \mathbf{l} | \mathbf{y}) = 0$
 1336

$$1337 S(\mathbf{x}; \mathbf{l}) + S(\mathbf{y}; \mathbf{l} | \mathbf{x}) = S(\mathbf{y}; \mathbf{l}) + S(\mathbf{x}; \mathbf{l} | \mathbf{y}),$$

1338
 1339 $S(\mathbf{y}; \mathbf{l} | \mathbf{x}) = 0 \implies S(\mathbf{x}; \mathbf{l}) = S(\mathbf{y}; \mathbf{l}) + S(\mathbf{x}; \mathbf{l} | \mathbf{y})$. $S(\mathbf{x}; \mathbf{l}) = 0 \implies -S(\mathbf{x}; \mathbf{l} | \mathbf{y}) = S(\mathbf{y}; \mathbf{l})$. Which
 1340 only holds if $S(\mathbf{x}; \mathbf{l} | \mathbf{y}) = S(\mathbf{y}; \mathbf{l}) = 0$.

1341 Similarly, if $S(\mathbf{x}; \mathbf{l} | \mathbf{y}) = 0$, then,
 1342

$$1343 S(\mathbf{y}; \mathbf{l} | \mathbf{x}) = -S(\mathbf{x}; \mathbf{l}) + S(\mathbf{y}; \mathbf{l})$$

1344 If $S(\mathbf{y}; \mathbf{l}) = 0$ then we are done. \square
 1345

1346 D.4 THEOREM 2
 1347

1348 **Theorem 2** (Spectral dependence). $\frac{1}{2} \text{Tr} [\mathbb{E} [\mathbb{E} [\Phi_1] : \Phi_1] - \mathbf{I}]$ is equal to $S(\mathbf{x}; \mathbf{l})$, if Φ is a canonical
 1349 Φ -embedding of the unconditional kind, or $S(\mathbf{y}; \mathbf{l} | \mathbf{x})$, if Φ is of the conditional kind.

1350 *Proof.* By corollary 8, $I_J(\mathbf{x}; \mathbf{I}) = \frac{1}{2} \mathbb{E} \left[D_J \left(\mathbb{P}(\mathbf{x} | \mathbf{I}) \parallel \mathbb{P}(\mathbf{x} | \mathbf{I}') \right) \right]$. Now, $D_J(\mathbb{P}_{\mathbf{x}|\mathbf{I}} \parallel \mathbb{P}_{\mathbf{x}|\mathbf{I}'}) =$
1351 $KL(\mathbb{P}_{\mathbf{x}|\mathbf{I}} \parallel \mathbb{P}_{\mathbf{x}|\mathbf{I}'}) + KL(\mathbb{P}_{\mathbf{x}|\mathbf{I}'} \parallel \mathbb{P}_{\mathbf{x}|\mathbf{I}})$. By proposition 12

$$1352 \quad KL(\gamma_{\mu_1, \Sigma_1} \parallel \gamma_{\mu_2, \Sigma_2}) = \frac{1}{2} \mathbf{D}(\Phi(\Sigma_1, \mu_1) \parallel \Phi(\Sigma_2, \mu_2)).$$

1353
1354 By definition 2 $\mathbf{D}(\Phi(\Sigma_1, \mu_1) \parallel \Phi(\Sigma_2, \mu_2)) = \mathbf{D}(\Phi_1 \parallel \Phi_{1'})$. Now, $\mathbf{D}(\Phi_1 \parallel \Phi_{1'}) + \mathbf{D}(\Phi_{1'} \parallel \Phi_1) =$
1355 $\text{Tr}[\Phi_1 : \Phi_{1'} + (\Phi_1 : \Phi_{1'})^{-1} - 2\mathbf{I}]$. By proposition 27, $\lambda(\Phi_1 : \Phi_{1'}) = \lambda(\Phi_{1'} : \Phi_1)$. Hence,

$$1356 \quad \mathbf{D}(\Phi_1 \parallel \Phi_{1'}) + \mathbf{D}(\Phi_{1'} \parallel \Phi_1) = \text{Tr}[\Phi_1 : \Phi_{1'} + \Phi_{1'} : \Phi_1 - 2\mathbf{I}].$$

1357 Thus $D_J(\mathbb{P}_{\mathbf{x}|\mathbf{I}} \parallel \mathbb{P}_{\mathbf{x}|\mathbf{I}'}) = \frac{1}{2} \text{Tr}[\Phi_1 : \Phi_{1'} + \Phi_{1'} : \Phi_1 - 2\mathbf{I}]$. Hence, $I_J(\mathbf{x}; \mathbf{I}) =$
1358 $\frac{1}{4} \mathbb{E} [\text{Tr}[\Phi_1 : \Phi_{1'} + \Phi_{1'} : \Phi_1 - 2\mathbf{I}]]$. By linearity of the trace and the expectation, $I_J(\mathbf{x}; \mathbf{I}) =$
1359 $\frac{1}{4} \text{Tr} [\mathbb{E} [\Phi_1 : \Phi_{1'}] + \mathbb{E} [\Phi_{1'} : \Phi_1] - 2\mathbf{I}]$. \mathbf{I} and \mathbf{I}' being independent. $\mathbb{E}[\Phi_1 : \Phi_{1'}] = \mathbb{E}[\Phi_{1'} : \Phi_1]$.
1360 Hence, $I_J(\mathbf{x}; \mathbf{I}) = \frac{1}{4} \text{Tr} [2\mathbb{E} [\Phi_1 : \Phi_{1'}] - 2\mathbf{I}]$ and $I_J(\mathbf{x}; \mathbf{I}) = \frac{1}{2} \text{Tr} [\mathbb{E} [\Phi_1 : \Phi_{1'}] - \mathbf{I}]$. This proves the
1361 statement for Φ -embeddings of the unconditional kind. The statement for Φ -embeddings of the
1362 conditional kind follow similarly by using proposition 13. \square

1363 D.4.1 SPECTRAL DIVERGENCES

1364 Relative entropies by looking at first differential of the logarithm of the size of a solution. The log
1365 determinant is concave by 28. We take the logarithm of the determinant to inject convexity in the
1366 scoring of perturbation along covariates and pairs of covariates, $\log \det [\mathbf{S}] = \text{Tr}[\mathbf{Log} \mathbf{S}]$. Taking the
1367 first differential, $d(\log \det [\mathbf{S}]; \mathbf{S}) = \text{Tr}[\mathbf{S}^{-1} d\mathbf{S}]$. Take $d\mathbf{S} \in \mathcal{S}_{\succ}(d)$, then $\text{Tr}[\mathbf{S}^{-1} d\mathbf{S}] = \text{Tr}[d\mathbf{S}\mathbf{S}^{-1}] =$
1368 $\text{Tr}[d\mathbf{S} : \mathbf{S}]$. Thus we have, $d(\log \det [\mathbf{S}]; d\mathbf{S}) = \text{Tr}[d\mathbf{S} : \mathbf{S}]$. Using the pairing, $\langle \mathbf{S}_1, \mathbf{S}_2 \rangle_F = \text{Tr}[\mathbf{S}_1 \mathbf{S}_2]$,
1369 the convex conjugate of the function $\mathbf{S} \mapsto F(\mathbf{S}) := -\log \det [\mathbf{S}]$, $\mathbf{Y} \mapsto \mathbf{F}^*(\mathbf{Y}) = -\log \det [\mathbf{I} - \mathbf{Y}]$.
1370 This allows establishing, conditions on the size of solution using Fenchel-Young type inequalities.
1371 Moreover, comparison of candidates solution is readily available through Bregman Divergences. The
1372 function F is precisely the negative Burg entropy (Burg, 1972). Turning \mathcal{Q} into a Bregman manifold.

1373 Log-det divergence

1374 **Proposition 11** (Burg's entropy Bregman's divergence). *Let $\mathbf{X} \mapsto F(\mathbf{X}) = -\log \det [\mathbf{X}]$, the*
1375 *\mathcal{Q} -negentropy generator, Let $\Phi_1, \Phi_2 \in \mathcal{Q}(d_1, d_2)$. Then, $B_F(\Phi_1, \Phi_2) = \text{Tr}[\Phi_1 : \Phi_2 - \mathbf{I}] -$*
1376 *$\log \det [\Phi_1 : \Phi_2]$.*

1377 *Proof.* By proposition 28 F is strictly convex, proper and lower semicontinuous on $\mathcal{S}_{\succeq}(d)$. Now
1378 $dF(\Phi_1, d\Phi_1) = -\text{Tr}[\Phi_1^{-1} d\Phi_1] = -\langle \Phi_1^{-1}, d\Phi_1 \rangle_F$. Hence $B_F(\Phi_1, \Phi_2) = -\log \det [\Phi_1] +$
1379 $\log \det [\Phi_2] + \langle \Phi_1 - \Phi_2, \Phi_2^{-1} \rangle_F = \log \det [\Phi_2 \Phi_1^{-1}] + \text{Tr}[(\Phi_1 - \Phi_2) \Phi_2^{-1}] = \text{Tr}[\Phi_2^{-1} \Phi_1 - \mathbf{I}] -$
1380 $\log \det [\Phi_2^{-1} \Phi_1]$. By proposition 26 $\lambda(\Phi_2^{-1} \Phi_1) = \lambda(\Phi_1 : \Phi_2)$. By proposition 6 $B_F(\Phi_1, \Phi_2) =$
1381 $\text{Tr}[\Phi_1 : \Phi_2 - \mathbf{I}] - \log \det [\Phi_1 : \Phi_2]$. \square

1382 *Remark 1.* The \mathcal{Q} -negentropy Bregman's divergence between Φ_1 and Φ_2 is the log-det divergence
1383 (Cichocki et al., 2015), $\mathbf{D}(\Phi_1 \parallel \Phi_2)$.

1384 This lead to the following spectral expressions,

1385 **Corollary 13.** *Let Λ be the eigenvalues matrix of $\Phi_1 : \Phi_2$. Then,*

$$1386 \quad \mathbf{D}(\Phi_1 \parallel \Phi_2) = \text{Tr}[\Lambda - \mathbf{I}] - \log \det [\Lambda] = \text{Tr}[\Lambda - \mathbf{Log} \Lambda - \mathbf{I}].$$

1387
1388 **Correspondance between log-det divergence Kullback-Leibler divergence** The next proposition
1389 shows that the KL-divergence of the distributions is precisely the log-det divergence of a Φ -embedding.
1390 **Proposition 12.**

$$1391 \quad KL(\gamma_{\mu_1, \Sigma_1} \parallel \gamma_{\mu_2, \Sigma_2}) = \frac{1}{2} \mathbf{D}(\Phi(\Sigma_1, \mu_1) \parallel \Phi(\Sigma_2, \mu_2)).$$

1404 *Proof.* By proposition 11 $\mathbf{D}(\Phi(\Sigma_1, \mu_1) \parallel \Phi(\Sigma_2, \mu_2)) = \text{Tr}[\Phi_1 : \Phi_2 - \mathbf{I}] - \log \det [\Phi_1 : \Phi_2]$. By
 1405 proposition 6, $\det [\Phi_1 : \Phi_2] = \det [\Sigma_2^{-1} \Sigma_1]$ and $\text{Tr}[\Phi_1 : \Phi_2] = \text{Tr}[\Sigma_2^{-1} \Sigma_1] + \text{Tr}[(\mu_2 - \mu_1)^\top \Sigma_2^{-1} (\mu_2 -$
 1406 $\mu_1)] + 1$. Hence $\mathbf{D}(\Phi_1 \parallel \Phi_2) = \log \left(\frac{\det[\Sigma_2]}{\det[\Sigma_1]} \right) + \text{Tr}[\Sigma_2^{-1} \Sigma_1] - m + \text{Tr}[(\mu_2 - \mu_1)^\top \Sigma_2^{-1} (\mu_2 - \mu_1)] =$
 1407 $2KL(\mathcal{N}(\mu_1, \Sigma_1) \parallel \mathcal{N}(\mu_2, \Sigma_2))$. \square
 1408
 1409

1410 **Symmetrized log-det divergence.** The log-det divergence is not symmetric, $\mathbf{D}(\Phi_1 \parallel \Phi_2) \neq$
 1411 $\mathbf{D}(\Phi_2 \parallel \Phi_1)$. However we can symmetrize it, by defining,
 1412

1413 **Definition 7** (J-divergence).

$$1414 \quad \mathbf{D}_J(\Phi_1 \parallel \Phi_2) = \mathbf{D}(\Phi_1 \parallel \Phi_2) + \mathbf{D}(\Phi_2 \parallel \Phi_1)$$

1415
 1416
 1417 **Proposition 13.**

$$1418 \quad \mathbf{D}_J(\Phi_1 \parallel \Phi_2) = \frac{1}{2} \text{Tr}[\Phi_1 : \Phi_2 + (\Phi_1 : \Phi_2)^{-1} - 2\mathbf{I}].$$

1419
 1420
 1421 *Proof.* Follows from definition 7, corollary 13 and proposition 26. \square
 1422

1423 Note that from Corollary 13 we get a spectral expression in terms of the eigenvalues $\Phi_1 : \Phi_2$.

1424 **Corollary 14** (J-Divergence spectral form).

$$1425 \quad D_J(\Phi_1 \parallel \Phi_2) = \frac{1}{2} \text{Tr}[\Lambda + \Lambda^{-1} - 2\mathbf{I}].$$

1426
 1427
 1428
 1429
 1430
 1431
 1432
 1433
 1434
 1435
 1436
 1437
 1438
 1439
 1440
 1441
 1442
 1443
 1444
 1445
 1446
 1447
 1448
 1449
 1450
 1451
 1452
 1453
 1454
 1455
 1456
 1457

1431 **Conditional Kullback-Leibler divergence.** It is easier to understand in the conditional Kullback-
 1432 Leilber as functional of two markov kernels and a probability distributions. The next proposition shows
 1433 that the conditional Kullback-Leibler divergence can also be expressed in terms of the log-determinant
 1434 divergence on \mathcal{Q} .

1435 **Proposition 14.** For $i \in \{1, 2\}$, Let \mathbf{K}_i a Gaussian Markov kernel with source \mathcal{X} , destination $\mathcal{P}(\mathcal{Y})$,
 1436 and $(\forall \mathbf{x} \in \mathcal{X}) \quad \mathbf{K}_i(\mathbf{x}) \sim \gamma_{\mu_i(\mathbf{x}), \Gamma_i}$, where $\mu_i(\mathbf{x}) := \mathbf{A}_i \mathbf{x} + \mathbf{b}_i$. Let Φ be the Φ -embedding of some
 1437 non-degenerate Gaussian measure $\gamma_{\mu, \Sigma}$. Define,

$$1438 \quad \Delta := [\mathbf{A}_2 - \mathbf{A}_1 \quad \mathbf{b}_2 - \mathbf{b}_1] \Phi \begin{bmatrix} (\mathbf{A}_2 - \mathbf{A}_1)^\top \\ (\mathbf{b}_2 - \mathbf{b}_1)^\top \end{bmatrix}.$$

1441
 1442
 1443
 1444
 1445
 1446
 1447
 1448
 1449
 1450
 1451
 1452
 1453
 1454
 1455
 1456
 1457

$$1442 \quad KL(\mathbf{K}_1 \parallel \mathbf{K}_2 \mid \gamma_{\mu, \Sigma}) = \frac{1}{2} \mathbf{D}(\Phi(\Gamma_1, \Delta) \parallel \Phi(\Gamma_2, \mathbf{0}_{d_y \times (d_x + 1)})).$$

1445
 1446
 1447
 1448
 1449
 1450
 1451
 1452
 1453
 1454
 1455
 1456
 1457

$$1447 \quad KL(\mathbf{K}_1 \parallel \mathbf{K}_2 \mid \gamma_{\mu, \Sigma}) = \mathbb{E} [KL(\mathbf{K}_1(\cdot, \mathbf{x}) \parallel \mathbf{K}_2(\cdot, \mathbf{x}))],$$

1449
 1450
 1451
 1452
 1453
 1454
 1455
 1456
 1457

$$1451 \quad (\forall \mathbf{x} \in \mathcal{X}) KL(\mathbf{K}_1(\cdot, \mathbf{x}) \parallel \mathbf{K}_2(\cdot, \mathbf{x})) = \frac{1}{2} \text{Tr} \left[\Phi_1(\mathbf{x}) : \Phi_2(\mathbf{x}) - (\mathbf{I}_{d_y} + 1) \right] - \log \det [\Phi_1(\mathbf{x}) : \Phi_2(\mathbf{x})].$$

1453
 1454
 1455
 1456
 1457

$$1454 \quad \text{Tr}[\Phi_1(\mathbf{x}) : \Phi_2(\mathbf{x})] = \text{Tr}[\Gamma_2^{-1} \Gamma_1] + \text{Tr}[(\mu_2(\mathbf{x}) - \mu_1(\mathbf{x}))^\top \Gamma_2^{-1} (\mu_2(\mathbf{x}) - \mu_1(\mathbf{x}))] + 1.$$

1456
 1457

$$1457 \quad \mathbb{E} [\text{Tr}[\Phi_1(\mathbf{x}) : \Phi_2(\mathbf{x})]] = \text{Tr}[\Gamma_2^{-1} \Gamma_1] + \text{Tr}[\Gamma_2^{-1} \mathbb{E} [(\mu_2(\mathbf{x}) - \mu_1(\mathbf{x}))(\mu_2(\mathbf{x}) - \mu_1(\mathbf{x}))^\top]] + 1. \quad (4)$$

1458 Writing $\mathbf{A} := \mathbf{A}_2 - \mathbf{A}_1$ and $\mathbf{b} := \mathbf{b}_2 - \mathbf{b}_1$. Expanding the quadratic product yields,

$$1459 \quad (\boldsymbol{\mu}_2(\mathbf{x}) - \boldsymbol{\mu}_1(\mathbf{x}))(\boldsymbol{\mu}_2(\mathbf{x}) - \boldsymbol{\mu}_1(\mathbf{x}))^\top = (\mathbf{A}\mathbf{x} + \mathbf{b})(\mathbf{A}\mathbf{x} + \mathbf{b})^\top,$$

$$1460 \quad = \mathbf{b}\mathbf{b}^\top + \mathbf{A}\mathbf{x}\mathbf{b}^\top + \mathbf{b}\mathbf{x}^\top\mathbf{A}^\top + \mathbf{A}\mathbf{x}\mathbf{x}^\top\mathbf{A}^\top.$$

1461 Taking the expectation we have,

$$1462 \quad \mathbb{E}[(\boldsymbol{\mu}_2(\mathbf{x}) - \boldsymbol{\mu}_1(\mathbf{x}))(\boldsymbol{\mu}_2(\mathbf{x}) - \boldsymbol{\mu}_1(\mathbf{x}))^\top] = \mathbf{b}\mathbf{b}^\top + \mathbf{A}\boldsymbol{\mu}\mathbf{A}^\top + \mathbf{A}\boldsymbol{\mu}\mathbf{b}^\top + \mathbf{A}\boldsymbol{\Sigma}\mathbf{A}^\top + \mathbf{b}\boldsymbol{\mu}^\top\mathbf{A}^\top,$$

$$1463 \quad = [\mathbf{A} \quad \mathbf{b}] \begin{bmatrix} \boldsymbol{\Sigma} + \boldsymbol{\mu}\boldsymbol{\mu}^\top & \boldsymbol{\mu} \\ \boldsymbol{\mu}^\top & 1 \end{bmatrix} \begin{bmatrix} \mathbf{A}^\top \\ \mathbf{b}^\top \end{bmatrix},$$

$$1464 \quad = [\mathbf{A} \quad \mathbf{b}] \boldsymbol{\Phi}(\boldsymbol{\Sigma}, \boldsymbol{\mu}) \begin{bmatrix} \mathbf{A}^\top \\ \mathbf{b}^\top \end{bmatrix}.$$

1465 Hence,

$$1466 \quad 2KL(\mathbf{K}_1 \parallel \mathbf{K}_2 \mid \gamma_{\boldsymbol{\mu}, \boldsymbol{\Sigma}}) = \text{Tr}[\boldsymbol{\Gamma}_2^{-1}\boldsymbol{\Gamma}_1 - \mathbf{I}_{d_y}] - \log \det [\boldsymbol{\Gamma}_2^{-1}\boldsymbol{\Gamma}_1] + \text{Tr} \left[\boldsymbol{\Gamma}_2^{-1} [\mathbf{A} \quad \mathbf{b}] \boldsymbol{\Phi}(\boldsymbol{\Sigma}, \boldsymbol{\mu}) \begin{bmatrix} \mathbf{A}^\top \\ \mathbf{b}^\top \end{bmatrix} \right].$$

1467 Writing, $\boldsymbol{\Delta} := [\mathbf{A} \quad \mathbf{b}] \boldsymbol{\Phi}(\boldsymbol{\Sigma}, \boldsymbol{\mu}) \begin{bmatrix} \mathbf{A}^\top \\ \mathbf{b}^\top \end{bmatrix}$, $\mathbf{M}_1 = \boldsymbol{\Phi}(\boldsymbol{\Gamma}_1, \boldsymbol{\Delta})$, $\mathbf{M}_2 = \boldsymbol{\Phi}(\boldsymbol{\Gamma}_2, \mathbf{0}_{d_y \times (d_x+1)})$ By proposition 6,

$$1468 \quad \text{Tr} [\boldsymbol{\Gamma}_2^{-1}\boldsymbol{\Delta}] = \text{Tr}[\mathbf{M}_1 : \mathbf{M}_2] - \text{Tr}[\boldsymbol{\Gamma}_2^{-1}\boldsymbol{\Gamma}_1] - (d_x + 1).$$

1469 Moreover, by proposition 6, $\det [\boldsymbol{\Gamma}_2^{-1}\boldsymbol{\Gamma}_1] = \det [\mathbf{M}_1 : \mathbf{M}_2]$ Thus we conclude,

$$1470 \quad KL(\mathbf{K}_1 \parallel \mathbf{K}_2 \mid \gamma_{\boldsymbol{\mu}, \boldsymbol{\Sigma}}) = \frac{1}{2} \mathbf{D} \left(\boldsymbol{\Phi}(\boldsymbol{\Gamma}_1, \boldsymbol{\Delta}) \parallel \boldsymbol{\Phi}(\boldsymbol{\Gamma}_2, \mathbf{0}_{d_y \times (d_x+1)}) \right).$$

1471 □

1482 D.5 COROLLARY 1

1483 **Corollary 1.** *Let $\boldsymbol{\Phi}_1$ be a state index canonical $\boldsymbol{\Phi}$ -embedding. Then,*

$$1484 \quad \mathbb{R}^{m+n} = \bigoplus_{\lambda \in \lambda(\mathbb{E}[\boldsymbol{\Phi}_1] : \boldsymbol{\Phi}_1)} \text{Eig} \left(\mathbb{E}[\mathbb{E}[\boldsymbol{\Phi}_1] : \boldsymbol{\Phi}_1], \lambda \right)$$

$$1485 \quad \min_{\mathbf{W}^\top \mathbf{W} = \mathbf{I}_k} \text{Tr} [\mathbf{W}^\top \mathbb{E}[\mathbb{E}[\boldsymbol{\Phi}_1] : \boldsymbol{\Phi}_1] \mathbf{W}] = \sum_{i=1}^k \lambda_i^\uparrow (\mathbb{E}[\mathbb{E}[\boldsymbol{\Phi}_1] : \boldsymbol{\Phi}_1]).$$

1486 *Proof.* Direct application of the Ky-Fan trace minimization principle (Fan, 1949). □

1487 D.6 THEOREM 3

1488 D.6.1 REDUCING STABILITY OF EIGENSPACES TO A MATRIX EQUATIONS

1489 How to find the condition under which a simple invariant space of of hermitian matrix is stable?
1490 Consider a symmetric matrix in block form,

$$1491 \quad \mathbf{S} = \begin{bmatrix} \mathbf{S}_{11} & \mathbf{S}_{12} \\ \mathbf{S}_{12}^\top & \mathbf{S}_{22} \end{bmatrix}.$$

1492 Let \mathcal{X} be a simple invariant subspace of \mathbf{S} . Write \mathbf{X}_1 for the matrix whose columns form an orthonormal basis fo \mathcal{X} and \mathbf{Y}_2 the matrix whose columns form an orthonormal basis for \mathcal{Y} . The matrix $\mathbf{W} = [\mathbf{X}_1 \quad \mathbf{Y}_2]$. \mathbf{W} is unitary. Indeed,

$$1493 \quad \mathbf{W}^\top \mathbf{W} = \begin{bmatrix} \mathbf{X}_1^\top \\ \mathbf{Y}_2^\top \end{bmatrix} [\mathbf{X}_1 \quad \mathbf{Y}_2] = \begin{bmatrix} \mathbf{X}_1^\top \mathbf{X}_1 & \mathbf{X}_1^\top \mathbf{Y}_2 \\ \mathbf{Y}_2^\top \mathbf{X}_1 & \mathbf{Y}_2^\top \mathbf{Y}_2 \end{bmatrix} = \begin{bmatrix} \mathbf{I} & \mathbf{0} \\ \mathbf{0} & \mathbf{I} \end{bmatrix}.$$

1494 Consider the *spectral resolution* of \mathbf{S} ,

$$1495 \quad \mathbf{W}^\top \mathbf{S} \mathbf{W} = \begin{bmatrix} \mathbf{X}_1^\top \mathbf{S} \mathbf{X}_1 & \mathbf{X}_1^\top \mathbf{S} \mathbf{Y}_2 \\ \mathbf{Y}_2^\top \mathbf{S} \mathbf{X}_1 & \mathbf{Y}_2^\top \mathbf{S} \mathbf{Y}_2 \end{bmatrix}.$$

1512 There exists diagonal matrices Λ_{11} and Λ_{22} such that $\mathbf{S}\mathbf{X}_1 = \mathbf{X}_1\Lambda_{11}$ and $\mathbf{S}\mathbf{Y}_2 = \mathbf{Y}_2\Lambda_{22}$. We therefore
 1513 have, $\mathbf{Y}_2^\top \mathbf{S}\mathbf{X}_1 = \mathbf{Y}_2^\top \mathbf{X}_1\Lambda_{11} = \mathbf{0}$, and hence,

$$1514 \mathbf{W}^\top \mathbf{S}\mathbf{W} = \begin{bmatrix} \Lambda_{11} & \mathbf{0} \\ \mathbf{0} & \Lambda_{22} \end{bmatrix}.$$

1517 Now consider a perturbation of $\tilde{\mathbf{S}} = \mathbf{S} + \mathbf{E}$. Partition \mathbf{E} conformably, and write,

$$1519 \mathbf{W}^\top \mathbf{E}\mathbf{W} = \begin{bmatrix} \mathbf{E}_{11} & \mathbf{E}_{12} \\ \mathbf{E}_{21} & \mathbf{E}_{22} \end{bmatrix}.$$

1521 we want to find conditions under which the invariant space \mathcal{X} of \mathbf{S} is also an invariant space of $\tilde{\mathbf{S}}$.
 1522 Form,

$$1523 \hat{\mathbf{X}}_1 = (\mathbf{X}_1 + \mathbf{Y}_2\mathbf{P})(\mathbf{I} + \mathbf{P}^\top\mathbf{P})^{-\frac{1}{2}},$$

$$1525 \hat{\mathbf{Y}}_2 = (\mathbf{Y}_2 - \mathbf{X}_1\mathbf{P}^\top)(\mathbf{I} + \mathbf{P}\mathbf{P}^\top)^{-\frac{1}{2}}.$$

1527 Write, $\hat{\mathbf{W}} = [\hat{\mathbf{X}}_1 \quad \hat{\mathbf{Y}}_2]$. Notice that $\hat{\mathbf{W}}$ is unitary. Indeed,

$$1528 \hat{\mathbf{W}}^\top \hat{\mathbf{W}} = \begin{bmatrix} (\mathbf{I} + \mathbf{P}^\top\mathbf{P})^{-1} + \sqrt{(\mathbf{I} + \mathbf{P}^\top\mathbf{P})^{-1}}\mathbf{P}^\top\mathbf{P}\sqrt{(\mathbf{I} + \mathbf{P}^\top\mathbf{P})^{-1}} & \mathbf{0} \\ \mathbf{0} & (\mathbf{I} + \mathbf{P}\mathbf{P}^\top)^{-1} + \sqrt{(\mathbf{I} + \mathbf{P}\mathbf{P}^\top)^{-1}}\mathbf{P}\mathbf{P}^\top\sqrt{(\mathbf{I} + \mathbf{P}\mathbf{P}^\top)^{-1}} \end{bmatrix}$$

$$1531 = \begin{bmatrix} \sqrt{(\mathbf{I} + \mathbf{P}^\top\mathbf{P})^{-1}}(\mathbf{I} + \mathbf{P}^\top\mathbf{P})\sqrt{(\mathbf{I} + \mathbf{P}\mathbf{P}^\top)^{-1}} & \mathbf{0} \\ \mathbf{0} & \sqrt{(\mathbf{I} + \mathbf{P}\mathbf{P}^\top)^{-1}}(\mathbf{I} + \mathbf{P}\mathbf{P}^\top)\sqrt{(\mathbf{I} + \mathbf{P}\mathbf{P}^\top)^{-1}} \end{bmatrix},$$

$$1534 = \begin{bmatrix} \mathbf{I} & \mathbf{0} \\ \mathbf{0} & \mathbf{I} \end{bmatrix}.$$

1536 Forming the resolution of $\tilde{\mathbf{S}}$ with respect to $\hat{\mathbf{W}}$, the lower left block is equal to zero if and only if,

$$1537 \mathbf{0} = - \left(\sqrt{(\mathbf{I} + \mathbf{P}\mathbf{P}^\top)^{-1}} (\mathbf{P}\Lambda_{11} + \mathbf{P}\mathbf{X}_1^\top\mathbf{E}\mathbf{X}_1 + \mathbf{P}\mathbf{X}_1^\top\mathbf{E}\mathbf{Y}_2\mathbf{P} - \mathbf{Y}_2^\top\mathbf{E}\mathbf{Y}_2\mathbf{P} - \mathbf{Y}_2^\top\mathbf{E}\mathbf{X}_1 - \Lambda_{22}\mathbf{P}) \sqrt{(\mathbf{I} + \mathbf{P}^\top\mathbf{P})^{-1}} \right), \Leftrightarrow$$

$$1539 \mathbf{0} = \mathbf{P}\Lambda_{11} + \mathbf{P}\mathbf{X}_1^\top\mathbf{E}\mathbf{X}_1 + \mathbf{P}\mathbf{X}_1^\top\mathbf{E}\mathbf{Y}_2\mathbf{P} - \mathbf{Y}_2^\top\mathbf{E}\mathbf{Y}_2\mathbf{P} - \mathbf{Y}_2^\top\mathbf{E}\mathbf{X}_1 - \Lambda_{22}\mathbf{0}.$$

$$1541 (\Lambda_{11} + \mathbf{X}_1^\top\mathbf{E}\mathbf{X}_1)\mathbf{P} - \mathbf{P}(\Lambda_{22} + \mathbf{Y}_2^\top\mathbf{E}\mathbf{Y}_2)\mathbf{P} = \mathbf{Y}_2\mathbf{E}\mathbf{X}_1 - \mathbf{P}\mathbf{X}_1^\top\mathbf{E}\mathbf{Y}_2\mathbf{P}. \quad (5)$$

1542 Let $\mathbf{L}_1 = \Lambda_{11} + \mathbf{E}_{11}$, $\mathbf{L}_2 = (\Lambda_{22} + \mathbf{E}_{22})$, $\mathbf{H} := \mathbf{E}_{12} = \mathbf{X}_1^\top\mathbf{E}\mathbf{Y}_2$, and $\mathbf{G} := \mathbf{E}_{21} = \mathbf{Y}_2^\top\mathbf{E}\mathbf{X}_1$. We can
 1543 rewrite equation (5) in terms of the Sylvester operator,

$$1544 \mathbf{T} : \mathbf{P} \mapsto \mathbf{P}\mathbf{L}_1 - \mathbf{L}_2\mathbf{P}, \quad (6)$$

1545 yielding,

$$1546 \mathbf{T}(\mathbf{P}) = \mathbf{G} - \mathbf{P}\mathbf{H}\mathbf{P}. \quad (7)$$

1548 D.6.2 EXISTENCE OF A SOLUTION

1549 **Proposition 15.** *The equation (7) admits a solution if,*

$$1551 4 \|\mathbf{G}\| \|\mathbf{H}\| \|\mathbf{T}^{-1}\|^2 < 1. \quad (8)$$

1554 *Proof.* Formally, Inverting \mathbf{T} in equation (7),

$$1556 \mathbf{P} = \mathbf{T}^{-1}(\mathbf{G} - \mathbf{P}\mathbf{H}\mathbf{P}), \quad (9)$$

1557 suggests that the solution to equation (7) should be a fixed point of equation (9). Form the sequence,

$$1558 \begin{cases} \mathbf{P}_0 & = \mathbf{0}, \\ \mathbf{P}_{k+1} & = \mathbf{T}^{-1}(\mathbf{G} - \mathbf{P}_k\mathbf{H}\mathbf{P}_k), \quad \forall k \in \mathbb{N} \cap [0, \infty). \end{cases} \quad (10)$$

1561 Consider set space of all matrices $\mathcal{M}_{m \times n}(\mathbb{R})$, endowed a metric derived from a unitarily invariant
 1562 norm,

$$1563 \mathcal{M}_{m \times n}(\mathbb{R}) \times \mathcal{M}_{m \times n}(\mathbb{R}) \rightarrow [0, \infty),$$

$$1564 (\mathbf{A}, \mathbf{B}) \mapsto d(\mathbf{A}, \mathbf{B}) := \|\mathbf{A} - \mathbf{B}\|.$$

1565 We need to establish when the sequence defined in equation (10):

1566

1. Cauchy.

1567

2. Converges to the fixed defined in equation (9).

1568

1569

First, we show that the sequence is bounded.

1570

$$\begin{aligned}
\|\mathbf{P}_{k+1}\| &\leq \|\mathbf{T}^{-1}\| \|\mathbf{G} - \mathbf{P}_k \mathbf{H} \mathbf{P}_k\|, \\
&\leq \|\mathbf{T}^{-1}\| (\|\mathbf{G}\| + \|\mathbf{P}_k \mathbf{H} \mathbf{P}_k\|), \\
&\leq \|\mathbf{T}^{-1}\| (\|\mathbf{G}\| + \|\mathbf{H}\| \|\mathbf{P}_k\|^2), \\
&= \|\mathbf{T}^{-1}\| \|\mathbf{G}\| + \|\mathbf{T}^{-1}\| \|\mathbf{H}\| \|\mathbf{P}_k\|^2
\end{aligned}$$

1576

Let $x_{k+1} = \|\mathbf{T}^{-1}\| \|\mathbf{G}\| + \|\mathbf{T}^{-1}\| \|\mathbf{H}\| x_k^2$ for $k \in \mathbb{N} \cap [0, \infty)$ and $x_0 = 0$. Then,

1577

$$\forall k \in \mathbb{N} \cap [0, \infty), x_{k+1} \geq x_k.$$

1578

Any accumulation point of x_k must verify,

1581

$$x = \|\mathbf{T}^{-1}\| \|\mathbf{G}\| + \|\mathbf{T}^{-1}\| \|\mathbf{H}\| x^2.$$

1582

Formally, The roots of this quadratic polynomial is x are given by,

1583

1584

$$x_{\pm} = \frac{1 \pm \sqrt{1 - 4 \|\mathbf{G}\| \|\mathbf{H}\| \|\mathbf{T}^{-1}\|^2}}{2 \|\mathbf{H}\| \|\mathbf{T}^{-1}\|}.$$

1585

1586

For the roots to be real we need to require,

1587

$$\Delta := 1 - 4 \|\mathbf{G}\| \|\mathbf{H}\| \|\mathbf{T}^{-1}\|^2 \geq 0 \quad (11)$$

1588

Since we are interested in a least upper bound, we will only consider the accumulation point x_- . Hence, we have,

1589

1590

$$\|\mathbf{P}\|_{k+1} \leq x_- = \frac{1 - \sqrt{1 - 4 \|\mathbf{G}\| \|\mathbf{H}\| \|\mathbf{T}^{-1}\|^2}}{2 \|\mathbf{H}\| \|\mathbf{T}^{-1}\|}. \quad (12)$$

1594

And hence, assuming $\Delta > 0$, then,

1595

$$\|\mathbf{P}\| \leq 2 \|\mathbf{G}\| \|\mathbf{T}^{-1}\|. \quad (13)$$

1596

The sequence in equation (10) is bounded and increasing.

1597

1598

$$\begin{aligned}
\|\mathbf{P}_{k+1} - \mathbf{P}_k\| &= \|\mathbf{T}^{-1}(\mathbf{G} - \mathbf{P}_k \mathbf{H} \mathbf{P}_k) - \mathbf{P}_k\|, \\
&= \|\mathbf{T}^{-1}(\mathbf{G} - \mathbf{P}_k \mathbf{H} \mathbf{P}_k) - (\mathbf{T}^{-1}(\mathbf{G} - \mathbf{P}_{k-1} \mathbf{H} \mathbf{P}_{k-1}))\|, \\
&\leq \|\mathbf{T}^{-1}\| \|\mathbf{P}_k \mathbf{H} \mathbf{P}_k - \mathbf{P}_{k-1} \mathbf{H} \mathbf{P}_{k-1}\|.
\end{aligned}$$

1599

1600

1601

1602

Now,

1603

$$\mathbf{P}_k \mathbf{H} \mathbf{P}_k - \mathbf{P}_{k-1} \mathbf{H} \mathbf{P}_{k-1} = \frac{1}{2} \{(\mathbf{P}_k - \mathbf{P}_{k-1}) \mathbf{H} (\mathbf{P}_k + \mathbf{P}_{k-1})\} + \frac{1}{2} \{(\mathbf{P}_k + \mathbf{P}_{k-1}) \mathbf{H} (\mathbf{P}_k - \mathbf{P}_{k-1})\}.$$

1604

1605

Taking the norm,

1606

1607

$$\begin{aligned}
\|\mathbf{P}_k \mathbf{H} \mathbf{P}_k - \mathbf{P}_{k-1} \mathbf{H} \mathbf{P}_{k-1}\| &= \left\| \frac{1}{2} ((\mathbf{P}_k - \mathbf{P}_{k-1}) \mathbf{H} (\mathbf{P}_k + \mathbf{P}_{k-1})) + \frac{1}{2} ((\mathbf{P}_k + \mathbf{P}_{k-1}) \mathbf{H} (\mathbf{P}_k - \mathbf{P}_{k-1})) \right\|, \\
&\leq \frac{1}{2} \left\{ \|(\mathbf{P}_k - \mathbf{P}_{k-1}) \mathbf{H} (\mathbf{P}_k + \mathbf{P}_{k-1})\| + \|(\mathbf{P}_k + \mathbf{P}_{k-1}) \mathbf{H} (\mathbf{P}_k - \mathbf{P}_{k-1})\| \right\}, \\
&\leq \frac{1}{2} \left\{ \|\mathbf{P}_k - \mathbf{P}_{k-1}\| \|\mathbf{H}\| \|\mathbf{P}_k + \mathbf{P}_{k-1}\| + \|\mathbf{P}_k + \mathbf{P}_{k-1}\| \|\mathbf{H}\| \|\mathbf{P}_k - \mathbf{P}_{k-1}\| \right\}, \\
&= \|\mathbf{H}\| \|\mathbf{P}_k + \mathbf{P}_{k-1}\| \|\mathbf{P}_k - \mathbf{P}_{k-1}\|, \\
&\leq 2 \max(\|\mathbf{P}_k\|, \|\mathbf{P}_{k-1}\|) \|\mathbf{H}\| \|\mathbf{P}_k - \mathbf{P}_{k-1}\|.
\end{aligned}$$

1608

1609

1610

1611

1612

Substituting, we have,

1613

1614

$$\begin{aligned}
\|\mathbf{P}_{k+1} - \mathbf{P}_k\| &\leq 2 \|\mathbf{T}^{-1}\| \max(\|\mathbf{P}_k\|, \|\mathbf{P}_{k-1}\|) \|\mathbf{H}\| \|\mathbf{P}_k - \mathbf{P}_{k-1}\|, \\
&\leq \left(2 \|\mathbf{T}^{-1}\| \|\mathbf{H}\| \frac{1 - \sqrt{\Delta}}{2 \|\mathbf{H}\| \|\mathbf{T}^{-1}\|} \right) \|\mathbf{P}_k - \mathbf{P}_{k-1}\|, \\
&= (1 - \sqrt{\Delta}) \|\mathbf{P}_k - \mathbf{P}_{k-1}\|.
\end{aligned}$$

1615

1616

1617

1618

1619

1620 Unpacking the recurrence,
1621

$$1622 \quad \|\mathbf{P}_{k+1} - \mathbf{P}_k\| \leq (1 - \sqrt{\Delta})^k \|\mathbf{P}_1 - \mathbf{P}_0\|.$$

1623
1624 Thus as long as $(1 - \sqrt{\Delta}) < 1$ which holds if and only if,
1625

$$1626 \quad 4 \|\mathbf{G}\| \|\mathbf{H}\| \|\mathbf{T}^{-1}\|^2 < 1, \quad (14)$$

1627
1628 the sequence in equation (10) is Cauchy. \square
1629

1630 D.6.3 THEOREM 3

1631 **Theorem 3** (Conditions for faithful processing). *Let \mathcal{E}_k be the rank- k subspace of minimal dependence*
1632 *of Φ . Let $\bar{\Phi} = \Phi + \mathbf{E}$. If $\|\mathbf{E}\|_F < \frac{1}{2} (\lambda_{k+1}^\uparrow(\Phi) - \lambda_k^\uparrow(\Phi))$, then \mathcal{E}_k is an algebraically invariant*
1633 *subspace of $\bar{\Phi}$.*
1634

1635
1636 *Proof.* Now that we have establish the conditions under which the matrix equation (7) we need to
1637 refine them. First let's express $\|\mathbf{T}^{-1}\|$. We re-express the sylvester operator. For any, \mathbf{P} we have,
1638

$$1639 \quad \begin{aligned} \mathbf{T}(\mathbf{P}) &= \mathbf{P}\mathbf{L}_1 - \mathbf{L}_2\mathbf{P}, \iff \\ \mathbf{vec}(\mathbf{T}(\mathbf{P})) &= \mathbf{vec}(\mathbf{P}\mathbf{L}_1 - \mathbf{L}_2\mathbf{P}), \\ &= (\mathbf{I} \otimes \mathbf{L}_1 - \mathbf{L}_2 \otimes \mathbf{I}) \mathbf{vec}(\mathbf{P}). \end{aligned}$$

1640
1641
1642 And hence, by propostion todo,
1643

$$1644 \quad \lambda_{ij}(\mathbf{I} \otimes \mathbf{L}_1 - \mathbf{L}_2^\top \otimes \mathbf{I}) = \lambda_i(\mathbf{L}_1) - \lambda_j(\mathbf{L}_2).$$

1645
1646 Indeed, by the spectral mapping theorem, (Kato, 1995), for any unitarily invariant norm, $\mathbf{X} \mapsto \|\mathbf{X}\|$
1647 there exists symmetric Gauge function g such that $\|\mathbf{X}\| = g(\sigma(\mathbf{X}))$.

$$1648 \quad \begin{aligned} \|\mathbf{T}^{-1}\| &= g(\sigma(\mathbf{T}^{-1})), \\ &= g(\sigma(\mathbf{T})^{-1}). \end{aligned}$$

1649
1650 Thus if consider the spectral norm then,
1651

$$1652 \quad \|\mathbf{T}^{-1}\| = \sigma_d(\mathbf{T}),$$

1653
1654 the seperation between \mathbf{L}_1 and \mathbf{L}_2 . Hence we have that,
1655

$$1656 \quad \text{sep}(\mathbf{L}_1, \mathbf{L}_2)^{-1} = \inf_{\mathbf{P} \neq \mathbf{0}} \frac{\|\mathbf{T}(\mathbf{P})\|_F}{\|\mathbf{P}\|_F},$$

1657
1658
1659 and hence,
1660

$$1661 \quad \text{sep}(\mathbf{L}_1, \mathbf{L}_2) = \min |\lambda(\mathbf{L}_{11}) - \lambda(\mathbf{L}_{22})|.$$

1662 Condition (14) becomes,
1663

$$1664 \quad 4 \|\mathbf{G}\| \|\mathbf{H}\| \text{sep}(\mathbf{L}_1, \mathbf{L}_2)^{-2} < 1. \quad (15)$$

1665 If we assume that $\mathbf{E} = \mathbf{E}^\top$ then $\mathbf{G}^\top = \mathbf{H}$ and hence, condition (15) becomes,
1666

$$1667 \quad 4 \|\mathbf{G}\|^2 \text{sep}(\mathbf{L}_1, \mathbf{L}_2)^{-2} < 1. \quad (16)$$

1668 or equivalently,
1669

$$1670 \quad 4 \left(\frac{\|\mathbf{E}_{12}\|}{\text{sep}(\mathbf{L}_1, \mathbf{L}_2)} \right)^2 < 1.$$

1671
1672 which holds if and only if,
1673

$$1674 \quad \|\mathbf{E}_{12}\| < \frac{\text{sep}(\mathbf{L}_1, \mathbf{L}_2)}{2}. \quad (17)$$

1674 Here I need to explicitly show that if condition (17) is satisfied then the spectra of the diagonal blocks
 1675 of the perturbed matrix are separated. In order to show the first point it is enough to establish that,
 1676

$$1677 \text{sep}(\mathbf{L}_1, \mathbf{L}_2) > 0.$$

1678 Indeed, we have,

$$1679 \begin{aligned} \mathbf{T}(\mathbf{P}) &= \mathbf{P}\mathbf{\Lambda}_{11} - \mathbf{\Lambda}_{22}\mathbf{P} + (\mathbf{P}\mathbf{X}_1^\top\mathbf{E}\mathbf{X}_1^\top - \mathbf{Y}_2^\top\mathbf{E}\mathbf{Y}_2\mathbf{P}), \\ 1680 \|\mathbf{T}(\mathbf{P})\| &\geq \|\mathbf{P}\mathbf{\Lambda}_{11} - \mathbf{\Lambda}_{22}\mathbf{P}\| - \|\mathbf{P}\mathbf{X}_1^\top\mathbf{E}\mathbf{X}_1^\top - \mathbf{Y}_2^\top\mathbf{E}\mathbf{Y}_2\mathbf{P}\|, \\ 1681 &\geq \|\mathbf{P}\mathbf{\Lambda}_{11} - \mathbf{\Lambda}_{22}\mathbf{P}\| - \|\mathbf{P}\mathbf{X}_1^\top\mathbf{E}\mathbf{X}_1\| - \|\mathbf{Y}_2^\top\mathbf{E}\mathbf{Y}_2\mathbf{P}\|, \\ 1682 &\geq \|\mathbf{P}\mathbf{\Lambda}_{11} - \mathbf{\Lambda}_{22}\mathbf{P}\| - \|\mathbf{P}\| \|\mathbf{X}_1^\top\mathbf{E}\mathbf{X}_1\| - \|\mathbf{P}\| \|\mathbf{Y}_2^\top\mathbf{E}\mathbf{Y}_2\|, \\ 1683 &\geq \|\mathbf{P}\mathbf{\Lambda}_{11} - \mathbf{\Lambda}_{22}\mathbf{P}\| - \|\mathbf{P}\| (\|\mathbf{X}_1^\top\mathbf{E}\mathbf{X}_1\| + \|\mathbf{Y}_2^\top\mathbf{E}\mathbf{Y}_2\|). \end{aligned}$$

1684 Thus taking the infimum over the set of all $\{\mathbf{P} : \|\mathbf{P}\| = 1\}$,

$$1685 \text{sep}(\mathbf{L}_1, \mathbf{L}_2) \geq \text{sep}(\mathbf{\Lambda}_{11}, \mathbf{\Lambda}_{22}) - (\|\mathbf{E}_{11}\| + \|\mathbf{E}_{22}\|).$$

1686 Moreover, by bounds (13) and (18), it is enough to require,

$$1687 \|\mathbf{E}\| < \frac{\text{sep}(\mathbf{\Lambda}_{11}, \mathbf{\Lambda}_{22})}{2}, \quad (18)$$

1688 If we choose the Froebenius norm, we have,

$$1689 \text{sep}(\mathbf{\Lambda}_{11}, \mathbf{\Lambda}_{22}) = \min|\lambda(\mathbf{\Lambda}_{11}) - \lambda(\mathbf{\Lambda}_{22})|.$$

1690 The condition reduces to,

$$1691 \|\mathbf{E}\| < \frac{\min|\lambda(\mathbf{\Lambda}_{11}) - \lambda(\mathbf{\Lambda}_{22})|}{2}.$$

1692 □

1703 D.7 PROPOSITION 2

1704 D.7.1 SIMPLIFICATION

1705 Two simplifications are helpful here. First, the analysis of **DGP-S** reduces to that of **DGP-U**. Second,
 1706 it is enough to consider only one ratio which will allow us to make the ratio more tractable by resorting
 1707 to the generalized eigenvalue problem.

1708 **Equivalence between DGP-U and DGP-S** Indeed, $y \sim \text{Rademecher}(\frac{1}{2})$, hence $\mathbb{E}[y] = 0$ and
 1709 $\mathbb{V}[y] = 1$. Hence, $\Phi_y = \mathbf{I}_2$. $\mathbf{A}_1 = \text{Cov}[\mathbf{x}, y | \mathbf{I}] \mathbb{V}[y]^{-1} = \mathbb{E}[\mathbf{x}y | \mathbf{I}] - \mathbb{E}[\mathbf{x} | \mathbf{I}] \mathbb{E}[y] = \mathbb{E}[\mathbf{x}y | \mathbf{I}]$. Now,

$$1710 \begin{aligned} \mathbb{E}[\mathbf{x}y | \mathbf{I}] &= \frac{1}{2} \mathbb{E}[\mathbf{x}y | \mathbf{I}, y = 1] + \frac{1}{2} \mathbb{E}[\mathbf{x}y | \mathbf{I}, y = -1], \\ 1711 &= \mathbf{A}^\top \boldsymbol{\mu} + \mathbf{B}^\top \boldsymbol{\mu}_1. \end{aligned}$$

1712 Hence, $\mathbf{A}_1 = \mathbf{A}^\top \boldsymbol{\mu} + \mathbf{B}^\top \boldsymbol{\mu}_1$. Following **DGP-0**, $\mathbf{x}_1 = \mathbf{A}_1 \mathbf{y} + \mathbf{b}_1 + \sqrt{\Gamma_1} \boldsymbol{\varepsilon}$. Taking the expectation we have,
 1713 $\mathbb{E}[\mathbf{x} | \mathbf{I}] = \mathbf{b}_1$, and thus,

$$1714 \begin{aligned} \mathbf{b}_1 &= \mathbb{E}[\mathbf{x} | \mathbf{I}], \\ 1715 &= \frac{1}{2} \mathbb{E}[\mathbf{A}^\top \boldsymbol{\mu} + \mathbf{B}^\top \boldsymbol{\mu}_1 | \mathbf{I}, y = 1] - \frac{1}{2} \mathbb{E}[\mathbf{A}^\top \boldsymbol{\mu} + \mathbf{B}^\top \boldsymbol{\mu}_1 | \mathbf{I}, y = -1], \\ 1716 &= \mathbf{0}. \end{aligned}$$

1717 Hence the arguments to the fundamental Φ -embedding is given by,

$$1718 \begin{aligned} \mathbf{X}_1 &= \mathbf{A}^\top \boldsymbol{\Sigma} \mathbf{A} + \mathbf{B}^\top \boldsymbol{\Sigma}_1 \mathbf{B}, \\ 1719 \mathbf{Y}_1 &= [(\mathbf{A}^\top \boldsymbol{\mu} + \mathbf{B}^\top \boldsymbol{\mu}_1) \mathbf{0}]. \end{aligned}$$

1720 In this setting the fundamental Φ -embedding of **DGP-U** and **DGP-S** coincide.

1728
1729
1730
1731
1732
1733
1734
1735
1736
1737
1738
1739
1740
1741
1742
1743
1744
1745
1746
1747
1748
1749
1750
1751
1752
1753
1754
1755
1756
1757
1758
1759
1760
1761
1762
1763
1764
1765
1766
1767
1768
1769
1770
1771
1772
1773
1774
1775
1776
1777
1778
1779
1780
1781

It is enough to consider a single ratio Consider the state dependence,

$$S = \frac{1}{2} \text{Tr}[\mathbb{E}[\mathbb{E}[\Phi_1] \bullet \Phi_1] - \mathbf{I}].$$

We have,

$$\begin{aligned} \text{Tr}[\mathbb{E}[\mathbb{E}[\Phi_1] : \Phi_1] - \mathbf{I}] &= \text{Tr}[\mathbb{E}[\Phi_1 : \Phi_1 - \mathbf{I}]], \\ &= \mathbb{P}(\mathbf{I} = 1)\mathbb{P}(\mathbf{I}' = 2) (\text{Tr}[\Phi_1 : \Phi_2 + \Phi_2 : \Phi_1 - 2\mathbf{I}]). \end{aligned}$$

By proposition 27 $\text{Tr}[\Phi_2 : \Phi_1] = \text{Tr}[(\Phi_1 : \Phi_2)^{-1}]$. Thus it is enough to look for a state independent subspace in $\Phi_1 : \Phi_2$. Diagonalizing (Φ_1, Φ_2) by congruence. There exists a non-singular matrix, Φ_2 -orthonormal, \mathbf{V} such that, $\Phi_2^{-1}\Phi_1 = \mathbf{V}\Lambda\mathbf{V}^{-1}$. Hence,

$$S = \frac{\mathbb{P}(\mathbf{I} = 1)\mathbb{P}(\mathbf{I} = 2)}{2} (\text{Tr}[\Lambda + \Lambda^{-1} - 2\mathbf{I}]).$$

Moreover, $\mathbf{W} = \Phi_2^{\frac{1}{2}}\mathbf{V}$ is orthonormal and Diagonalizes $\Phi_1 : \Phi_2$.

D.7.2 TWO USEFUL LEMMAS

Lemma 4.

$$\mathbf{A}\mathbf{A}^\top - \mathbf{B}\mathbf{B}^\top = \frac{1}{2}(\mathbf{A} - \mathbf{B})(\mathbf{A}^\top + \mathbf{B}^\top) + \frac{1}{2}(\mathbf{A} + \mathbf{B})(\mathbf{A}^\top - \mathbf{B}^\top).$$

Proof. Direct computation. □

Under **DGP-U**,

Lemma 5.

$$\Phi_{x|i} = (\mathbf{F} \oplus 1)^\top \Phi_{z|i} (\mathbf{F} \oplus 1).$$

Proof. Direct computation. □

D.7.3 DGP-U ADMITS 1 AS AN EIGENVALUE

The next proposition shows that for any two different states i, j the state conditional pencils of Φ -embeddings under **DGP-U** admit 1 as an eigenvalue with geometric multiplicity equal to d_c .

Proposition 16 (1 as eigenvalue for DGP-U). *Let $i \neq j \in |\mathcal{L}|$. Then, $1 \in \lambda(\Phi_{x|i=i}, \Phi_{x|i=j})$ and $m(1) = d_c$.*

Proof. First we have, $\Phi = (\mathbf{F} \oplus 1)^\top (\Phi_i - \lambda\Phi_j)(\mathbf{F} \oplus 1)$, $\mathbf{F} = \begin{bmatrix} \mathbf{A} \\ \mathbf{B} \end{bmatrix} \in \text{Hom}(\mathbb{R}^d, \mathbb{R}^{d_c+d_i})$. (Rosenfeld et al., 2020) requires \mathbf{F}^\top to be injective; $\dim(\text{Ker}(\mathbf{F}^\top)) = 0$ and subsequently also surjective. Hence, the injectivity assumption on \mathbf{F} requires,

$$d = d_c + d_i. \tag{19}$$

Moreover, $\text{rank } \mathbf{F} = \text{rank } \mathbf{A} + \text{rank } \mathbf{B} - \dim(\mathcal{R}(\mathbf{A}^\top) \cap \mathcal{R}(\mathbf{B}^\top))$. We therefore $\text{rank } \mathbf{F} = d_c + d_i - \dim(\mathcal{R}(\mathbf{A}^\top) \cap \mathcal{R}(\mathbf{B}^\top))$. By equation (19), $\dim(\mathcal{R}(\mathbf{A}^\top) \cap \mathcal{R}(\mathbf{B}^\top)) = 0$.

Form the characteristic polynomial for $(\Phi_{x|i=i}, \Phi_{x|i=j})$, $p(\lambda) := \det[\Phi_{x|i=i} - \lambda\Phi_{x|i=j}]$. Let $\Phi_i = \Phi_{z|i=i}$, from lemma 5, $\Phi_{x|i=i} - \lambda\Phi_{x|i=j} = (\mathbf{F} \oplus 1)^\top (\Phi_i - \lambda\Phi_j)(\mathbf{F} \oplus 1)$. Expanding Φ -embedding for Φ_i ,

$$\Phi_i = \begin{bmatrix} \mu_c \mu_c^\top + \Sigma_c & \mu_c \mu_i^\top & \mu_c \\ \mu_i \mu_c^\top & \mu_i \mu_i^\top + \Sigma_i & \mu_i \\ \mu_c^\top & \mu_i^\top & 1 \end{bmatrix}.$$

1782 This yields,

$$1783 \Phi_i - \lambda \Phi_j = \begin{bmatrix} (1 - \lambda)(\mu_c \mu_c^\top) + \Sigma_c - \lambda \Sigma_c & \mu_c(\mu_i^\top - \lambda \mu_j^\top) & \mu_c - \lambda \mu_c \\ 1784 (\mu_i - \lambda \mu_j) \mu_c^\top & \mu_i \mu_i^\top + \Sigma_i - \lambda \Sigma_j - \lambda(\mu_j \mu_j^\top) & \mu_i - \lambda \mu_j \\ 1785 \mu_c^\top - \lambda \mu_c^\top & \mu_i^\top - \lambda \mu_j^\top & 1 - \lambda \end{bmatrix}.$$

1786 In particular if $\lambda = 1$,

$$1787 \Phi_i - \Phi_j = \begin{bmatrix} \mathbf{0} & \mu_c(\mu_i^\top - \mu_j^\top) & \mathbf{0} \\ 1788 (\mu_i - \mu_j) \mu_c^\top & \mu_i \mu_i^\top + \Sigma_i - \Sigma_j - \mu_j \mu_j^\top & \mu_i - \mu_j \\ 1789 \mathbf{0} & \mu_i^\top - \mu_j^\top & \mathbf{0} \end{bmatrix}. \quad (20)$$

1790 We have $\Phi := \Phi_{x||=i} - \Phi_{x||=j} = (\mathbf{F} \oplus 1)^\top (\Phi_i - \Phi_j) (\mathbf{F} \oplus 1)$. Expanding,

$$1791 \Phi = \begin{bmatrix} (\mathbf{B}^\top (\mu_i \mu_i^\top + \Sigma_i - \Sigma_j - \mu_j \mu_j^\top) + \mathbf{A}^\top \mu_c (\mu_i^\top - \mu_j^\top)) \mathbf{B} + \mathbf{B}^\top (\mu_i - \mu_j) \mu_c^\top \mathbf{A} & \mathbf{B}^\top (\mu_i - \mu_j) \\ 1792 (\mu_i^\top - \mu_j^\top) \mathbf{B} & \mathbf{0} \end{bmatrix}.$$

1793 By lemma 5

$$1794 \Phi = (\mathbf{B} \oplus 1)^\top \left[\begin{bmatrix} ((\mu_i \mu_i^\top + \Sigma_i - \Sigma_j - \mu_j \mu_j^\top) + \mathbf{A}^\top \mu_c (\mu_i^\top - \mu_j^\top)) + (\mu_i - \mu_j) \mu_c^\top \mathbf{A} & (\mu_i - \mu_j) \\ 1795 (\mu_i^\top - \mu_j^\top) & \mathbf{0} \end{bmatrix} (\mathbf{B} \oplus 1) \right].$$

1800 Now $\mathbf{B} \in \mathcal{M}_{d_l \times d}(\mathbb{R})$, by equation (19) $d_l < d$. Hence, $\text{rank } \mathbf{B} \leq d_l < d$ and $\text{rank}(\mathbf{B} \oplus 1) = \text{rank } \mathbf{B} + 1 \leq$
 1801 $d_l + 1 < d + 1$. Thus $\text{rank}(\Phi_{x||=i} - \Phi_{x||=j}) \leq d_l + 1$ which implies $p(1) = \det[\Phi_{x||=i} - \Phi_{x||=j}] = 0$
 1802 and thus $1 \in \lambda(\Phi_{x||=i}, \Phi_{x||=j})$.

1803 The (geometric) multiplicity of the eigenvalue 1 of Φ is equal to the dimensionality of $\text{Ker}(\Phi)$. Φ
 1804 being finite dimensional it is enough to look at its rank.

1805 Now, by the injectivity assumption $\text{rank}(\mathbf{F} \oplus 1) = \text{rank}(\mathbf{F}) + 1$. $\mathbf{F} \oplus 1$ being non-singular; $\text{rank}(\Phi) =$
 1806 $\text{rank}(\Phi_i - \Phi_j)$. It is enough to consider the rank of $\Phi_i - \Phi_j$.

$$1807 \Phi_i - \Phi_j = \begin{bmatrix} \mathbf{0} & \mu_c(\mu_i^\top - \mu_j^\top) & \mathbf{0} \\ 1808 (\mu_i - \mu_j) \mu_c^\top & \mu_i \mu_i^\top + \Sigma_i - \Sigma_j - \mu_j \mu_j^\top & \mu_i - \mu_j \\ 1809 \mathbf{0} & \mu_i^\top - \mu_j^\top & \mathbf{0} \end{bmatrix}.$$

1810 By block Gauss Jordan elimination $\Phi_i - \Phi_j$ is congruent to,

$$1811 \begin{bmatrix} (\mu_i - \mu_j) \mu_c^\top & \mu_i \mu_i^\top + \Sigma_i - \Sigma_j - \mu_j \mu_j^\top & \mu_i - \mu_j \\ 1812 \mathbf{0} & \mu_i^\top - \mu_j^\top & \mathbf{0} \\ 1813 \mathbf{0} & \mathbf{0} & \mathbf{0} \end{bmatrix}.$$

1814 The rank of $\Phi_i - \Phi_j$ can then be read directly and hence $\text{rank}(\Phi_i - \Phi_j) = d_l + 1 = \text{rank}(\Phi)$. By the
 1815 fundamental theorem of linear maps, $\dim \text{Ker}(\Phi) = d_c$. \square

1820 D.7.4 RESTRICTION TO INDEPENDENT SUBSPACE

1821 Next we show that the restriction of the Φ -embeddings to \mathcal{E}_1 are equal.

1822 **Proposition 17.** *Let $(\Phi_{x||=i}, \Phi_{x||=j})$ be a pencil of Φ -embeddings under DGP-U. Let \mathcal{E}_1 the pencils*
 1823 *eigenspace associated to $\lambda = 1$. Then,*

$$1824 \Phi_{x||=i}|_{\mathcal{E}_1} = \Phi_{x||=j}|_{\mathcal{E}_1}.$$

1825 *Proof.* proposition 16 establishes that $(\Phi_{x||=i}, \Phi_{x||=j})$ admits an eigenspace \mathcal{E}_1 . Let
 1826 $(\Lambda = \mathbf{I} \oplus \Lambda_2, \mathbf{V} = [\mathbf{V}_1 \ \mathbf{V}_2])$ be an eigensystem of $(\Phi_{x||=i}, \Phi_{x||=j})$. First the matrix \mathbf{V} is non-singular.
 1827 Let's assume that it is $\Phi_{x||=j}$ normalized. Consider the spectral resolutions,

$$1828 \begin{bmatrix} \mathbf{V}_1^\top \\ \mathbf{V}_2^\top \end{bmatrix} \Phi_{x||=i} [\mathbf{V}_1 \ \mathbf{V}_2] = \begin{bmatrix} \mathbf{I} & \mathbf{0} \\ \mathbf{0} & \Lambda_2 \end{bmatrix},$$

$$1829 \begin{bmatrix} \mathbf{V}_1^\top \\ \mathbf{V}_2^\top \end{bmatrix} \Phi_{x||=j} [\mathbf{V}_1 \ \mathbf{V}_2] = \begin{bmatrix} \mathbf{I} & \mathbf{0} \\ \mathbf{0} & \mathbf{I} \end{bmatrix}.$$

From here we will deduce that $\Phi_{x|i=1}|_{\mathcal{E}_1} = \Phi_{x|i=j}|_{\mathcal{E}_1}$. Let's express the restrictions of the Φ -embeddings to \mathcal{E}_1 . First we need an injection from \mathbb{R}^{m+n} to \mathcal{E}_1 . $\mathbf{V}_1 \in \text{Hom}(\mathbb{R}^{d_c}, \mathbb{R}^d)$. Hence, $J_{\mathcal{E}_1} = \mathbf{V}_1$. The restrictions to \mathcal{E}_1 , $\Phi_{x|i=1}|_{\mathcal{E}_1} = \mathbf{V}_1^\dagger \Phi_{x|i=1} \mathbf{V}_1 = \Phi_{x|i=j}|_{\mathcal{E}_1} = \mathbf{V}_1^\dagger \Phi_{x|i=j} \mathbf{V}_1$. We have

$$\Phi_{x|i=1}|_{\mathcal{E}_1} = \Phi_{x|i=j}|_{\mathcal{E}_1} \iff \mathbf{V}_1^\dagger \Phi_{x|i=1} \mathbf{V}_1 = \mathbf{V}_1^\dagger \Phi_{x|i=j} \mathbf{V}_1.$$

Noting that the columns of \mathbf{V}_1 being independent, $\mathbf{V}_1^\dagger = (\mathbf{V}_1^\top \mathbf{V}_1)^{-1} \mathbf{V}_1^\top$. Hence,

$$\Phi_{x|i=1}|_{\mathcal{E}_1} = \Phi_{x|i=j}|_{\mathcal{E}_1} \iff (\mathbf{V}_1^\top \mathbf{V}_1)^{-1} \mathbf{V}_1^\top \Phi_{x|i=1} \mathbf{V}_1 = (\mathbf{V}_1^\top \mathbf{V}_1)^{-1} \mathbf{V}_1^\top \Phi_{x|i=j} \mathbf{V}_1 \iff \mathbf{V}_1^\top (\Phi_{x|i=1} - \Phi_{x|i=j}) \mathbf{V}_1 = \mathbf{0}.$$

□

D.7.5 CONDITION ON DGP-U MOMENTS

Now that we have an equivalent condition for the equality of the restriction to \mathcal{E}_1 . We will exploit the block spectral structure of pencils of Φ -embeddings in order to deduce from equality of the restrictions of the Φ -embedding on \mathcal{E}_1 constraints on the form of **DGP-U**.

Proposition 18. *Let $(\Phi_{x|i=i}, \Phi_{x|i=j})$ be a pencil of Φ -embeddings under **DGP-U**. Let \mathcal{E}_1 the pencils eigenspace associated to $\lambda = 1$. Let \mathbf{V}_1 be a $\Phi_{x|i=j}$ -orthonormal matrix such that $\mathcal{R}(\mathbf{V}_1) = \mathcal{E}_1$.*

Partition \mathbf{V}_1 conformably, $\mathbf{V}_1 = \begin{bmatrix} \mathbf{V}_{11} \\ \mathbf{V}_{21} \end{bmatrix}$. Then,

$$\begin{cases} \mathbf{V}_{11}^\top \mathbf{B}^\top (\boldsymbol{\Sigma}_i - \boldsymbol{\Sigma}_j) \mathbf{B} \mathbf{V}_{11} & = \mathbf{0}, \\ \mathbf{V}_{11}^\top \mathbf{B}^\top (\boldsymbol{\mu}_i - \boldsymbol{\mu}_j) & = \mathbf{0}. \end{cases} \quad (21)$$

Proof. Partition \mathbf{V}_1 conformably, $\mathbf{V}_1 = \begin{bmatrix} \mathbf{V}_{11} \\ \mathbf{V}_{21} \end{bmatrix}$. We have $\Phi_{x|i=k} = \Phi(\mathbf{A}_k, \mathbf{B}_k)$ for some $\mathbf{A}_k \in \mathcal{S}_>(m)$ and $\mathbf{B}_k \in \mathbb{R}^{m \times n}$. By proposition 17, $\mathbf{V}_1^\top (\Phi(\mathbf{A}_i, \mathbf{B}_i) - \Phi(\mathbf{A}_j, \mathbf{B}_j)) \mathbf{V}_1 = \mathbf{0}$, if and only,

$$\begin{aligned} \mathbf{V}_{11}^\top (\mathbf{B}_i \mathbf{B}_i^\top + \mathbf{A}_i - \mathbf{A}_j - \mathbf{B}_j \mathbf{B}_j^\top) \mathbf{V}_{11} + \mathbf{V}_{21}^\top (\mathbf{B}_i^\top - \mathbf{B}_j^\top) \mathbf{V}_{11} + \mathbf{V}_{11}^\top (\mathbf{B}_i - \mathbf{B}_j) \mathbf{V}_{21} &= \mathbf{0}, \\ \mathbf{V}_{11}^\top (\mathbf{A}_i - \mathbf{A}_j) \mathbf{V}_{11} + \mathbf{V}_{11}^\top (\mathbf{B}_i \mathbf{B}_i^\top - \mathbf{B}_j \mathbf{B}_j^\top) \mathbf{V}_{11} + \mathbf{V}_{21}^\top (\mathbf{B}_i^\top - \mathbf{B}_j^\top) \mathbf{V}_{11} + \mathbf{V}_{11}^\top (\mathbf{B}_i - \mathbf{B}_j) \mathbf{V}_{21} &= \mathbf{0}. \end{aligned}$$

By lemma 4,

$$\mathbf{B}_i \mathbf{B}_i^\top - \mathbf{B}_j \mathbf{B}_j^\top = \frac{1}{2} (\mathbf{B}_i - \mathbf{B}_j) (\mathbf{B}_i^\top + \mathbf{B}_j^\top) + \frac{1}{2} (\mathbf{B}_i + \mathbf{B}_j) (\mathbf{B}_i^\top - \mathbf{B}_j^\top).$$

Hence, $\mathbf{V}_1^\top (\Phi(\mathbf{A}_i, \mathbf{B}_i) - \Phi(\mathbf{A}_j, \mathbf{B}_j)) \mathbf{V}_1 = \mathbf{0}$, if and only,

$$\begin{aligned} &\mathbf{V}_{11}^\top (\mathbf{A}_i - \mathbf{A}_j) \mathbf{V}_{11} + \\ &\mathbf{V}_{11}^\top \left(\frac{1}{2} (\mathbf{B}_i - \mathbf{B}_j) (\mathbf{B}_i^\top + \mathbf{B}_j^\top) + \frac{1}{2} (\mathbf{B}_i + \mathbf{B}_j) (\mathbf{B}_i^\top - \mathbf{B}_j^\top) \right) \mathbf{V}_{11} + \\ &\mathbf{V}_{21}^\top (\mathbf{B}_i^\top - \mathbf{B}_j^\top) \mathbf{V}_{11} + \\ &\mathbf{V}_{11}^\top (\mathbf{B}_i - \mathbf{B}_j) \mathbf{V}_{21} = \mathbf{0}. \end{aligned}$$

By proposition 18, $\mathbf{V}_{11}^\top (\mathbf{B}_i - \mathbf{B}_j) = \mathbf{0}$. Hence, $\mathbf{V}_1^\top (\Phi(\mathbf{A}_i, \mathbf{B}_i) - \Phi(\mathbf{A}_j, \mathbf{B}_j)) \mathbf{V}_1 = \mathbf{0}$, if and only,

$$\begin{cases} \mathbf{V}_{11}^\top (\mathbf{A}_i - \mathbf{A}_j) \mathbf{V}_{11} & = \mathbf{0}, \\ \mathbf{V}_{11}^\top (\mathbf{B}_i - \mathbf{B}_j) & = \mathbf{0}. \end{cases} \quad (22)$$

Now in our setting, from **DGP-U** $\mathbf{A}_k = \mathbf{A}^\top \boldsymbol{\Sigma}_c \mathbf{A} + \mathbf{B}^\top \boldsymbol{\Sigma}_k \mathbf{B}$, and $\mathbf{B}_k = \mathbf{A}^\top \boldsymbol{\mu}_c + \mathbf{B}^\top \boldsymbol{\mu}_k$. Substituting into equation (22), we conclude

$$\begin{cases} \mathbf{V}_{11}^\top \mathbf{B}^\top (\boldsymbol{\Sigma}_i - \boldsymbol{\Sigma}_j) \mathbf{B} \mathbf{V}_{11} & = \mathbf{0}, \\ \mathbf{V}_{11}^\top \mathbf{B}^\top (\boldsymbol{\mu}_i - \boldsymbol{\mu}_j) & = \mathbf{0}. \end{cases} \quad (23)$$

□

1890 D.7.6 RECOVERY MECHANISM

1891 In this section we show explicitly how the spurious parts of **DGP-U** is annihilated by the eigenvectors
 1892 spanning \mathcal{E}_\perp . Then we show in which sense is the common/invariant parts of **DGP-U** is recovered.

1893 **Proposition 19.** Under **DGP-U**, $\mathbf{B}\mathbf{V}_{11} = \mathbf{0}$.

1894 *Proof.* By equation (21),

$$1895 \mathbf{V}_{11}^\top \mathbf{B}^\top (\boldsymbol{\Sigma}_i - \boldsymbol{\Sigma}_j) \mathbf{B} \mathbf{V}_{11} = \mathbf{0}. \quad (24)$$

1896 By lemma 6 there exists non-singular matrix \mathbf{Q} diagonalizing $(\boldsymbol{\Sigma}_i, \boldsymbol{\Sigma}_j)$. Take \mathbf{Q} to be $\boldsymbol{\Sigma}_j$ -
 1897 orthonormal. Then $\mathbf{Q}^\top \boldsymbol{\Sigma}_i \mathbf{Q} = \mathbf{D}$ and $\mathbf{Q}^\top \boldsymbol{\Sigma}_j \mathbf{Q} = \mathbf{I}$. Substituting in equation (24),

$$1901 \mathbf{V}_{11}^\top \mathbf{B}^\top \mathbf{Q}^{-\top} (\mathbf{D} - \mathbf{I}) \mathbf{Q}^{-1} \mathbf{B} \mathbf{V}_{11} = \mathbf{0}. \quad (25)$$

1902 Let $\mathbf{M}^\top := \mathbf{Q}^{-1} \mathbf{B} \mathbf{V}_{11}$ and equation (25) reads,

$$1903 \mathbf{M} \mathbf{D} \mathbf{M}^\top = \mathbf{M} \mathbf{M}^\top. \quad (26)$$

1904 Let see equation (26) as linear matrix equation in \mathbf{D} . By lemma 7 equation (26) admits a solution if
 1905 and only if,

$$1906 \mathbf{M} \mathbf{M}^\dagger \mathbf{M} \mathbf{M}^\top \mathbf{M}^{\dagger\top} \mathbf{M}^\top = \mathbf{M} \mathbf{M}^\top, \Leftrightarrow$$

$$1907 \mathbf{M} \mathbf{M}^\dagger \mathbf{M} \mathbf{M}^\top = \mathbf{M} \mathbf{M}^\top, \Leftrightarrow$$

$$1908 \mathbf{M} \mathbf{M}^\top = \mathbf{M} \mathbf{M}^\top.$$

1909 The general solution is therefore given by,

$$1910 \mathbf{D} = \mathbf{M}^\dagger \mathbf{M} \mathbf{M}^\top \mathbf{M}^{\dagger\top} + \mathbf{Q} - \mathbf{M}^\dagger \mathbf{M} \mathbf{Q} \mathbf{M}^\top \mathbf{M}^{\dagger\top},$$

$$1911 \mathbf{X} = \mathbf{M}^\dagger \mathbf{M} + \mathbf{Q} - \mathbf{M}^\dagger \mathbf{M} \mathbf{Q} \mathbf{M}^\top \mathbf{M}^{\dagger\top}.$$

1912 For any arbitrary matrix \mathbf{Q} . In particular for $\mathbf{Q} = \mathbf{0}$, the solution $\mathbf{X} = \mathbf{M}^\dagger \mathbf{M}$ is the orthogonal
 1913 projection onto $\mathcal{R}(\mathbf{M}^\top)$. It is now enough to show that \mathbf{D} cannot satisfy this requirement. Indeed,
 1914 for \mathbf{D} to be an orthogonal projection it must verify $\forall i \in \{1, \dots, d_l\} \mathbf{D}_{ii} \in \{0, 1\}$. But by assumption,
 1915 \mathbf{D} has at least one element different than 1. Moreover, since $\boldsymbol{\Sigma}_1$ and $\boldsymbol{\Sigma}_2$ are assumed to be positive
 1916 definite the generalized eigenvalues of the pencil are strictly greater than zero. Thus, \mathbf{D} cannot be a
 1917 solution equation (26); The equation therefore requires that $\mathbf{M}^\top = \mathbf{0}$. \mathbf{Q} being non-singular this is
 1918 only possible this means that $\mathbf{B}\mathbf{V}_{11} = \mathbf{0}$. \square

1919 **Proposition 20.** Under **DGP-U**. Let \mathbf{V}_1 be the eigenmatrix spanning the eigenspace \mathcal{E}_\perp of
 1920 $(\Phi_{\mathbf{x}|i}, \Phi_{\mathbf{x}|j})$. Then,

$$1921 \mathcal{R}(\mathbf{V}_{11}) = \mathcal{R}(\mathbf{A}^\top).$$

1922 *Proof.* By proposition 4 $\Phi_{\mathbf{x}|i} > \mathbf{0}$ and hence $(\Phi_{\mathbf{x}|i}/\mathbf{I}) > \mathbf{0}$. Moreover, $\mathbb{V}[\mathbf{x} | \mathbf{I}] > \mathbf{0}$. Hence,
 1923 $\text{Ker}(\mathbb{V}[\mathbf{x} | \mathbf{I} = i]) \cap \text{Ker}(\mathbb{V}[\mathbf{x} | \mathbf{I} = j]) = \{\mathbf{0}\}$. Under **DGP-U**, $\mathbb{V}[\mathbf{x} | \mathbf{I} = i] = \mathbf{A}^\top \boldsymbol{\Sigma}_c \mathbf{A} + \mathbf{B}^\top \boldsymbol{\Sigma}_i \mathbf{B}$.
 1924 Diagonalize $\mathbf{A}^\top \boldsymbol{\Sigma}_c \mathbf{A}$ and $\mathbf{B}^\top \boldsymbol{\Sigma}_i \mathbf{B}$ by congruence. There there exists a non-singular matrix \mathbf{X} such
 1925 that,

$$1926 \mathbf{X}^\top \mathbb{V}[\mathbf{x} | \mathbf{I} = i] \mathbf{X} = \mathbf{Diag}(\boldsymbol{\alpha}) + \mathbf{Diag}(\boldsymbol{\beta}).$$

1927 By Sylvester's theorem of inertia, $\boldsymbol{\alpha}$ has exactly $\text{rank}(\mathbf{A}^\top \boldsymbol{\Sigma}_c \mathbf{A}) = d_c$ positive entries and d_s zero entries.
 1928 Similarly $\mathbf{B}^\top \boldsymbol{\Sigma}_i \mathbf{B}$ has exactly $\text{rank}(\mathbf{B}^\top \boldsymbol{\Sigma}_i \mathbf{B}) = d_s$ positive entries and d_c zero entries. $\mathbf{X}^\top \mathbb{V}[\mathbf{x} | \mathbf{I} = i] \mathbf{X}$
 1929 is also diagonal and by positive definiteness of $\mathbb{V}[\mathbf{x} | \mathbf{I} = i]$ all of its entries are positive. By the
 1930 pigeon-hole principle, there is no single index $i \in [d_c + d_s]$ such that $\alpha_i = \beta_i = 0$. Hence,

$$1931 \text{Ker}(\mathbf{A}^\top \boldsymbol{\Sigma}_c \mathbf{A}) \cap \text{Ker}(\mathbf{B}^\top \boldsymbol{\Sigma}_i \mathbf{B}) \cap \text{Ker}(\mathbf{B}^\top \boldsymbol{\Sigma}_j \mathbf{B}) = \{\mathbf{0}\}. \quad (27)$$

1932 Let \mathbf{v} be any eigenvector of \mathcal{E}_\perp . proposition 19 implies that $\mathbf{v} \in \text{Ker}(\mathbf{B})$ and hence $\mathbf{v} \in \text{Ker}(\mathbf{B}^\top \boldsymbol{\Sigma}_j \mathbf{B})$.
 1933 By equation (27), $\mathbf{v} \notin \text{Ker}(\mathbf{A}^\top \boldsymbol{\Sigma}_c \mathbf{A})$; $\mathbf{v} \in \text{Ker}(\mathbf{A}^\top \boldsymbol{\Sigma}_c \mathbf{A})^\perp = \mathcal{R}(\mathbf{A}^\top \boldsymbol{\Sigma}_c \mathbf{A})$. $\boldsymbol{\Sigma}_c$ being positive definite,

1944 $\mathcal{R}(\mathbf{A}^\top \boldsymbol{\Sigma}_c \mathbf{A}) = \mathcal{R}\left(\mathbf{A}^\top \boldsymbol{\Sigma}_c^{\frac{1}{2}} \boldsymbol{\Sigma}_c^{\frac{1}{2}} \mathbf{A}\right)$. Now, $\mathcal{R}(\mathbf{A}^\top \boldsymbol{\Sigma}_c^{\frac{1}{2}} \boldsymbol{\Sigma}_c^{\frac{1}{2}} \mathbf{A}) = \mathcal{R}\left(\mathbf{A}^\top \boldsymbol{\Sigma}_c^{\frac{1}{2}}\right)$. Hence, $\mathbf{v} \in \mathcal{R}(\mathbf{A}^\top \boldsymbol{\Sigma}_c^{\frac{1}{2}})$.
1945
1946 Therefore, Writing \mathbf{V}_{11} be which columns are d_c eigenvectors of \mathcal{E}_1 , $\mathcal{R}(\mathbf{V}_{11}) \subseteq \mathcal{R}(\mathbf{A}^\top \boldsymbol{\Sigma}_c^{\frac{1}{2}})$. By
1947 proposition 16, $\text{rank}(\mathcal{R}(\mathbf{V}_{11})) = d_c = \text{rank}\left(\boldsymbol{\Sigma}_c^{\frac{1}{2}}\right) = \text{rank}(\mathbf{A} \boldsymbol{\Sigma}_c^{\frac{1}{2}})$. Thus, $\mathcal{R}(\mathbf{V}_{11}) = \mathcal{R}\left(\mathbf{A}^\top \boldsymbol{\Sigma}_c^{\frac{1}{2}}\right) =$
1948
1949 $\mathbf{A}^\top \mathcal{R}(\boldsymbol{\Sigma}_c^{\frac{1}{2}}) = \mathbf{A}^\top \mathbb{R}^{d_c} = \mathcal{R}(\mathbf{A}^\top)$. \square
1950
1951

1952 D.7.7 PROPOSITION 2

1953
1954 It is now enough to take $\mathbf{P} = \mathbf{V}_{11} \mathbf{V}_{11}^\dagger$ the orthogonal projection onto $\mathcal{R}(\mathbf{V}_{11})$ and we conclude,
1955 **Proposition 2.** *The expected ratio $\mathbb{E}[\mathbb{E}[\Phi_1] : \Phi_1]$ of both **DGP-U** and **DGP-S** admit a state independent subspace \mathcal{E}_\perp of dimension d_c . Moreover, a matrix $\mathbf{W}_1^\top = [\mathbf{W}_{11}^\top \quad \mathbf{W}_{12}^\top]$ such that $\mathcal{R}(\mathbf{W}_{11}) = \mathcal{E}_\perp$. Then the orthogonal projection \mathbf{P} onto $\mathcal{R}(\mathbf{W}_{11})$ verifies,*

- 1959 (i) $\mathbf{B}\mathbf{P} = \mathbf{0}$ (Annihilation of state dependent component),
- 1960 (ii) $\mathcal{R}(\mathbf{A}^\top) = \mathcal{R}(\mathbf{P})$ (Preservation of oracle component),
- 1961 (iii) $\mathbf{P}\mathbf{x} = \mathbf{A}^\top \mathbf{c}$ for **DGP-U**, and, $\mathbf{P}\mathbf{x} = \mathbf{y}\mathbf{A}^\top \mathbf{c}$ for **DGP-S**.

1963 E MATHEMATICAL ELEMENTS

1964 E.1 SCHUR COMPLEMENT

1965
1966 **Definition 8.** Let α and β multi-indices of rows and columns of a given matrix $\mathbf{M} \in \mathcal{M}_{r \times c}(\mathbb{R})$. Let α^c, β^c denote their set complement in the row and column indices in $[r]$ and $[c]$ respectively. Define the generalized Schur-complement the matrix $\mathbf{M}[\alpha, \beta]$ in \mathbf{M} is

$$1967 \mathbf{M}/\mathbf{M}[\alpha, \beta] = \mathbf{M}[\alpha^c, \beta^c] - \mathbf{M}[\alpha^c, \beta] \mathbf{M}[\alpha, \beta]^\dagger \mathbf{M}[\alpha, \beta^c].$$

1972 When $\alpha = \beta$, the β multi-index is omitted, and the generalized schur-complement reads,

$$1973 \mathbf{M}/\mathbf{M}[\alpha] := \mathbf{M}[\alpha^c, \alpha^c] - \mathbf{M}[\alpha^c, \alpha] \mathbf{M}[\alpha, \alpha]^\dagger \mathbf{M}[\alpha, \alpha^c].$$

1974 E.2 LDU AND UDL DECOMPOSITION

1975
1976 **Proposition 21** (UDL/LDU decompositions).

$$1977 \begin{bmatrix} \mathbf{I} & \mathbf{0} \\ \mathbf{C}\mathbf{A}^{-1} & \mathbf{I} \end{bmatrix} \begin{bmatrix} \mathbf{A} & \mathbf{0} \\ \mathbf{0} & \mathbf{D} - \mathbf{C}\mathbf{A}^{-1}\mathbf{B} \end{bmatrix} \begin{bmatrix} \mathbf{I} & \mathbf{A}^{-1}\mathbf{B} \\ \mathbf{0} & \mathbf{I} \end{bmatrix} = \begin{bmatrix} \mathbf{A} & \mathbf{B} \\ \mathbf{C} & \mathbf{D} \end{bmatrix} = \begin{bmatrix} \mathbf{I} & \mathbf{B}\mathbf{D}^{-1} \\ \mathbf{0} & \mathbf{I} \end{bmatrix} \begin{bmatrix} \mathbf{A} - \mathbf{B}\mathbf{D}^{-1}\mathbf{C} & \mathbf{0} \\ \mathbf{0} & \mathbf{D} \end{bmatrix} \begin{bmatrix} \mathbf{I} & \mathbf{0} \\ \mathbf{D}^{-1}\mathbf{C} & \mathbf{I} \end{bmatrix}$$

1983
1984 *Proof.* Let $\mathbf{X} := \begin{bmatrix} \mathbf{A} & \mathbf{B} \\ \mathbf{C} & \mathbf{D} \end{bmatrix}$. First form the pivot Multiply the first row by \mathbf{A}^{-1} , $\mathbf{E}_1 = \begin{bmatrix} \mathbf{A}^{-1} & \mathbf{0} \\ \mathbf{0} & \mathbf{I} \end{bmatrix}$.

$$1985 \mathbf{E}_1 \mathbf{X} = \begin{bmatrix} \mathbf{I} & \mathbf{A}^{-1}\mathbf{B} \\ \mathbf{C} & \mathbf{D} \end{bmatrix}.$$

1986
1987 Subtract \mathbf{C} times the first row from the second, $\mathbf{E}_2 = \begin{bmatrix} \mathbf{I} & \mathbf{0} \\ -\mathbf{C} & \mathbf{I} \end{bmatrix}$.

$$1988 \mathbf{E}_2 \mathbf{E}_1 \mathbf{X} = \begin{bmatrix} \mathbf{A} & \mathbf{B} \\ \mathbf{0} & \mathbf{D} - \mathbf{C}\mathbf{A}^{-1}\mathbf{B} \end{bmatrix}.$$

1989
1990 Multiply the first row by \mathbf{A} , $\mathbf{E}_3 = \begin{bmatrix} \mathbf{A} & \mathbf{0} \\ \mathbf{0} & \mathbf{I} \end{bmatrix}$.

$$1991 \mathbf{E}_3 \mathbf{E}_2 \mathbf{E}_1 \mathbf{X} = \begin{bmatrix} \mathbf{A} & \mathbf{B} \\ \mathbf{0} & \mathbf{D} - \mathbf{C}\mathbf{A}^{-1}\mathbf{B} \end{bmatrix}.$$

1998
1999
2000
2001
2002
2003
2004
2005
2006
2007
2008
2009
2010
2011
2012
2013
2014
2015
2016
2017
2018
2019
2020
2021
2022
2023
2024
2025
2026
2027
2028
2029
2030
2031
2032
2033
2034
2035
2036
2037
2038
2039
2040
2041
2042
2043
2044
2045
2046
2047
2048
2049
2050
2051

Now proceed similarly from the right and multiply the first column by \mathbf{A}^{-1} , $\mathbf{F}_1 = \begin{bmatrix} \mathbf{I} & \mathbf{A}^{-1} \\ \mathbf{0} & \mathbf{I} \end{bmatrix}$,

$$\mathbf{E}_3 \mathbf{E}_2 \mathbf{E}_1 \mathbf{X} \mathbf{F}_1 = \begin{bmatrix} \mathbf{I} & \mathbf{B} \\ \mathbf{0} & \mathbf{D} - \mathbf{C} \mathbf{A}^{-1} \mathbf{B} \end{bmatrix}.$$

Then the second column gets $-\mathbf{B}$ the first, $\mathbf{F}_2 = \begin{bmatrix} \mathbf{I} & \mathbf{0} \\ -\mathbf{B} & \mathbf{I} \end{bmatrix}$,

$$\mathbf{E}_3 \mathbf{E}_2 \mathbf{E}_1 \mathbf{X} \mathbf{F}_1 \mathbf{F}_2 = \begin{bmatrix} \mathbf{I} & \mathbf{0} \\ \mathbf{0} & \mathbf{D} - \mathbf{C} \mathbf{A}^{-1} \mathbf{B} \end{bmatrix}.$$

Finally the first column get multiplied by \mathbf{A} , $\mathbf{F}_3 = \begin{bmatrix} \mathbf{A} & \mathbf{0} \\ \mathbf{0} & \mathbf{I} \end{bmatrix}$.

$$\mathbf{E}_3 \mathbf{E}_2 \mathbf{E}_1 \mathbf{X} \mathbf{F}_1 \mathbf{F}_2 \mathbf{F}_3 = \begin{bmatrix} \mathbf{A} & \mathbf{0} \\ \mathbf{0} & \mathbf{D} - \mathbf{C} \mathbf{A}^{-1} \mathbf{B} \end{bmatrix}.$$

Thus we have,

$$\mathbf{E} := \mathbf{E}_1 \mathbf{E}_2 \mathbf{E}_3 = \begin{bmatrix} \mathbf{I} & \mathbf{0} \\ -\mathbf{C} \mathbf{A}^{-1} & \mathbf{I} \end{bmatrix}, \quad \mathbf{F} := \mathbf{F}_1 \mathbf{F}_2 \mathbf{F}_3 = \begin{bmatrix} \mathbf{I} & -\mathbf{A}^{-1} \mathbf{B} \\ \mathbf{0} & \mathbf{I} \end{bmatrix}.$$

Moreover, $\mathbf{E}^{-1} = \begin{bmatrix} \mathbf{I} & \mathbf{0} \\ \mathbf{C} \mathbf{A}^{-1} & \mathbf{I} \end{bmatrix}$ and $\mathbf{F}^{-1} = \begin{bmatrix} \mathbf{I} & \mathbf{A}^{-1} \mathbf{B} \\ \mathbf{0} & \mathbf{I} \end{bmatrix}$. Hence,

$$\begin{bmatrix} \mathbf{A} & \mathbf{B} \\ \mathbf{C} & \mathbf{D} \end{bmatrix} = \begin{bmatrix} \mathbf{I} & \mathbf{0} \\ \mathbf{C} \mathbf{A}^{-1} & \mathbf{I} \end{bmatrix} \begin{bmatrix} \mathbf{A} & \mathbf{0} \\ \mathbf{0} & \mathbf{D} - \mathbf{C} \mathbf{A}^{-1} \mathbf{B} \end{bmatrix} \begin{bmatrix} \mathbf{I} & \mathbf{A}^{-1} \mathbf{B} \\ \mathbf{0} & \mathbf{I} \end{bmatrix}.$$

To get a UDL decomposition start by using the second row as a pivot, $\mathbf{E}_1 = \begin{bmatrix} \mathbf{I} & \mathbf{0} \\ \mathbf{0} & \mathbf{D}^{-1} \end{bmatrix}$,

$$\mathbf{E}_1 \mathbf{X} = \begin{bmatrix} \mathbf{A} & \mathbf{B} \\ \mathbf{D}^{-1} \mathbf{C} & \mathbf{I} \end{bmatrix}.$$

Then the first row gets $-\mathbf{B}$ the second, $\mathbf{E}_2 = \begin{bmatrix} \mathbf{I} & -\mathbf{B} \\ \mathbf{0} & \mathbf{I} \end{bmatrix}$.

$$\mathbf{E}_2 \mathbf{E}_1 \mathbf{X} = \begin{bmatrix} \mathbf{A} - \mathbf{B} \mathbf{D}^{-1} \mathbf{C} & \mathbf{0} \\ \mathbf{D}^{-1} \mathbf{C} & \mathbf{I} \end{bmatrix}$$

The second row gets multiplied by \mathbf{D} , $\mathbf{E}_3 = \begin{bmatrix} \mathbf{I} & \mathbf{0} \\ \mathbf{0} & \mathbf{D} \end{bmatrix}$.

$$\mathbf{E}_3 \mathbf{E}_2 \mathbf{E}_1 \mathbf{X} = \begin{bmatrix} \mathbf{A} - \mathbf{B} \mathbf{D}^{-1} \mathbf{C} & \mathbf{0} \\ \mathbf{C} & \mathbf{D} \end{bmatrix}.$$

Next we start operating on the columns. The second column gets multiplied by \mathbf{D}^{-1} , $\mathbf{F}_1 = \begin{bmatrix} \mathbf{I} & \mathbf{0} \\ \mathbf{0} & \mathbf{D}^{-1} \end{bmatrix}$.

$$\mathbf{E}_3 \mathbf{E}_2 \mathbf{E}_1 \mathbf{X} \mathbf{F}_1 = \begin{bmatrix} \mathbf{A} - \mathbf{B} \mathbf{D}^{-1} \mathbf{C} & \mathbf{0} \\ \mathbf{C} & \mathbf{I} \end{bmatrix}.$$

The first column gets $-\mathbf{C}$ the second, $\mathbf{F}_2 = \begin{bmatrix} \mathbf{I} & \mathbf{0} \\ -\mathbf{C} & \mathbf{I} \end{bmatrix}$.

$$\mathbf{E}_3 \mathbf{E}_2 \mathbf{E}_1 \mathbf{X} \mathbf{F}_1 \mathbf{F}_2 = \begin{bmatrix} \mathbf{A} - \mathbf{B} \mathbf{D}^{-1} \mathbf{C} & \mathbf{0} \\ \mathbf{0} & \mathbf{I} \end{bmatrix}.$$

2052 Finally the second column gets multiplied by \mathbf{D} , $\mathbf{F}_3 = \begin{bmatrix} \mathbf{I} & \mathbf{0} \\ \mathbf{0} & \mathbf{D} \end{bmatrix}$.

$$2055 \mathbf{E}_3 \mathbf{E}_2 \mathbf{E}_1 \mathbf{X} \mathbf{F}_1 \mathbf{F}_2 \mathbf{F}_3 = \begin{bmatrix} \mathbf{A} - \mathbf{B} \mathbf{D}^{-1} \mathbf{C} & \mathbf{0} \\ \mathbf{0} & \mathbf{D} \end{bmatrix}.$$

$$2059 \mathbf{E} := \mathbf{E}_1 \mathbf{E}_2 \mathbf{E}_3 = \begin{bmatrix} \mathbf{I} & -\mathbf{B} \mathbf{D} \\ \mathbf{0} & \mathbf{I} \end{bmatrix}, \quad \mathbf{F} := \mathbf{F}_1 \mathbf{F}_2 \mathbf{F}_3 = \begin{bmatrix} \mathbf{I} & \mathbf{0} \\ -\mathbf{D}^{-1} \mathbf{C} & \mathbf{I} \end{bmatrix}.$$

2061 Hence,

$$2062 \begin{bmatrix} \mathbf{A} & \mathbf{B} \\ \mathbf{C} & \mathbf{D} \end{bmatrix} = \begin{bmatrix} \mathbf{I} & \mathbf{B} \mathbf{D}^{-1} \\ \mathbf{0} & \mathbf{I} \end{bmatrix} \begin{bmatrix} \mathbf{A} - \mathbf{B} \mathbf{D}^{-1} \mathbf{C} & \mathbf{0} \\ \mathbf{0} & \mathbf{D} \end{bmatrix} \begin{bmatrix} \mathbf{I} & \mathbf{0} \\ \mathbf{D}^{-1} \mathbf{C} & \mathbf{I} \end{bmatrix}.$$

2066 \square

2067 E.3 SIMULTANEOUS DIAGONALIZATION OF SPD MATRICES BY CONGRUENCE

2068 **Lemma 6.** *Let $\mathbf{S}_1, \mathbf{S}_2 \in \mathcal{S}_{++}(d)$. Then there exists a matrix \mathbf{V} diagonalizing both \mathbf{S}_1 and \mathbf{S}_2 by*
 2069 *congruence while also diagonalizing their quotient, $\mathbf{S}_2^{-1} \mathbf{S}_1$, by similarity.*

2071 *Proof.* By the spectral theorem, $\mathbf{D}_1^{-1/2} \mathbf{V}_1^\top \mathbf{S}_1 \mathbf{V}_1 \mathbf{D}_1^{-1/2} = \mathbf{I}$ and $\mathbf{D}_1^{-1/2} \mathbf{V}_1^\top \mathbf{S}_2 \mathbf{V}_1 \mathbf{D}_1^{-1/2} = \mathbf{V} \mathbf{A} \mathbf{V}^\top$. The matrix
 2072 $\mathbf{V}_1 \mathbf{A}_1^{-1/2} \mathbf{V}$ simultaneously diagonalize \mathbf{S}_1 and \mathbf{S}_2 by congruence. \square

2075 E.4 SPECTRAL EQUIVALENCES

2076 **Proposition 22** (Eigenpairs of the inverse). *Let $\mathbf{A} \in \mathcal{S}_{>}$. Then,*

- 2077 (i) \mathbf{A} and \mathbf{A}^{-1} share the same set of eigenvectors,
- 2078 (ii) $\lambda(\mathbf{A}^{-1}) = \lambda(\mathbf{A})^{-1}$, and,
- 2079 (iii) $\lambda^\downarrow(\mathbf{A}^{-1}) = \lambda^\uparrow(\mathbf{A})^{-1}$.

2081 **Proposition 23.** *Let $\mathbf{A}, \mathbf{B} \in \mathcal{S}_{>}(d)$. Then,*

- 2082 (i) If \mathbf{V} is an eigenmatrix of (\mathbf{A}, \mathbf{B}) then $\mathbf{B}^{1/2} \mathbf{V}$ is one for $\mathbf{B}^{-1/2} \mathbf{A} \mathbf{B}^{-1/2}$,
- 2083 (ii) $\lambda(\mathbf{A}, \mathbf{B}) = \lambda(\mathbf{B}^{-1/2} \mathbf{A} \mathbf{B}^{-1/2})$, and,
- 2084 (iii) $\lambda^\downarrow(\mathbf{A}, \mathbf{B}) = \lambda^\downarrow(\mathbf{B}^{-1/2} \mathbf{A} \mathbf{B}^{-1/2})$.

2085 *Proof.* \mathbf{B} being positive definite it has a unique non singular square root. Let (λ, \mathbf{v}) be an eigenpair of
 2086 (\mathbf{A}, \mathbf{B}) if and only, $\mathbf{A} \mathbf{v} = \lambda \mathbf{B} \mathbf{v} \iff \mathbf{B}^{-1/2} \mathbf{A} \mathbf{v} = \lambda \mathbf{B}^{1/2} \mathbf{v}$ (iii). Writing $\mathbf{w} = \mathbf{B}^{1/2} \mathbf{v}$, then (λ, \mathbf{v}) if and only
 2087 if $\mathbf{B}^{-1/2} \mathbf{A} \mathbf{B}^{-1/2} \mathbf{w} = \lambda \mathbf{w}$, which holds if and only if (λ, \mathbf{w}) is an eigenpair of $\mathbf{B}^{-1/2} \mathbf{A} \mathbf{B}^{-1/2}$, this proves (i),
 2088 (ii), and (iii). \square

2089 **Proposition 24.** *Let $\mathbf{A}, \mathbf{B} \in \mathcal{S}_{>}(d)$. Then,*

- 2090 (i) (\mathbf{A}, \mathbf{B}) and (\mathbf{B}, \mathbf{A}) share the same eigenvectors,
- 2091 (ii) $\lambda(\mathbf{B}, \mathbf{A}) = \lambda(\mathbf{A}, \mathbf{B})^{-1}$,
- 2092 (iii) $\lambda^\downarrow(\mathbf{B}, \mathbf{A}) = \lambda^\uparrow(\mathbf{A}, \mathbf{B})$.

2093 *Proof.* By proposition 23 (λ, \mathbf{v}) is an eigenpair of (\mathbf{A}, \mathbf{B}) if and only if it verifies $\mathbf{A} \mathbf{v} = \lambda \mathbf{B} \mathbf{v}$, if and
 2094 only if, $\mathbf{B} \mathbf{v} = \frac{1}{\lambda} \mathbf{A} \mathbf{v}$ if and only if $(\frac{1}{\lambda}, \mathbf{v})$ is an eigenpair of (\mathbf{B}, \mathbf{A}) (i)(ii). The map $t \mapsto \frac{1}{t}$ being order
 2095 reversing, $\lambda^\downarrow(\mathbf{B}, \mathbf{A}) = \lambda^\uparrow(\mathbf{A}, \mathbf{B})$. \square

2096 **Proposition 25.** *Let $\mathbf{A}, \mathbf{B} \in \mathcal{S}_{>}(d)$. Let \mathbf{L} be block lower triangular and \mathbf{D} be block diagonal with*
 2097 *positive definite blocks such that $\mathbf{B} = \mathbf{L}^\top \mathbf{D} \mathbf{L}$. Then,*

- 2098 (i) If \mathbf{V} is an eigenmatrix of (\mathbf{A}, \mathbf{B}) then $\mathbf{D}^{1/2} \mathbf{L} \mathbf{V}$ is one for $\mathbf{D}^{-1/2} \mathbf{L}^{-1} \mathbf{A} \mathbf{L}^{-\top} \mathbf{D}^{-1/2}$,

2106
2107
2108
2109
2110
2111
2112
2113
2114
2115
2116
2117
2118
2119
2120
2121
2122
2123
2124
2125
2126
2127
2128
2129
2130
2131
2132
2133
2134
2135
2136
2137
2138
2139
2140
2141
2142
2143
2144
2145
2146
2147
2148
2149
2150
2151
2152
2153
2154
2155
2156
2157
2158
2159

- (ii) $\lambda(\mathbf{A}, \mathbf{B}) = \lambda(\mathbf{D}^{-\frac{1}{2}} \mathbf{L}^{-1} \mathbf{A} \mathbf{L}^{-\top} \mathbf{D}^{-\frac{1}{2}})$, and,
 (iii) $\lambda^\downarrow(\mathbf{A}, \mathbf{B}) = \lambda^\uparrow(\mathbf{L}^{-1} \mathbf{D}^{-\frac{1}{2}} \mathbf{A} \mathbf{L}^{-\top} \mathbf{D}^{-\frac{1}{2}})$.

Proof. Write, $\mathbf{C} := \mathbf{L}^\top \mathbf{D}^{\frac{1}{2}}$ then $\mathbf{B} = \mathbf{C} \mathbf{C}^\top$ (λ, \mathbf{v}) is an eigenpair of (\mathbf{A}, \mathbf{B}) if and only, $\mathbf{A} \mathbf{v} = \lambda \mathbf{B} \mathbf{v} \iff \mathbf{C}^{-1} \mathbf{A} \mathbf{v} = \lambda \mathbf{C}^\top \mathbf{v}$ (iii). Writing $\mathbf{w} := \mathbf{C}^\top \mathbf{v}$, (λ, \mathbf{v}) is an eigenpair of (\mathbf{A}, \mathbf{B}) if and only if $\mathbf{C}^{-1} \mathbf{A} \mathbf{C}^{-\top} \mathbf{w} = \lambda \mathbf{w}$, which holds if and only if (λ, \mathbf{w}) is an eigenpair of $\mathbf{C}^{-1} \mathbf{A} \mathbf{C}^{-\top}$, this proves (i), (ii), and (iii). \square

Proposition 26 (Spectral equivalence $:$). *Let $\Phi_1, \Phi_2 \in \mathcal{Q}$. Then,*

- (i) $\lambda(\Phi_1 : \Phi_2) = \lambda(\Phi_1, \Phi_2)$, and,
 (ii) $\lambda^\downarrow(\Phi_1 : \Phi_2) = \lambda^\downarrow(\Phi_1, \Phi_2)$.

Proof. (i) and (ii) follow by the definition of $:$ and proposition 25. \square

Proposition 27. *Let $\Phi_1, \Phi_2 \in \mathcal{Q}$. Then,*

- (i) $\lambda(\Phi_2 : \Phi_1) = \lambda(\Phi_1 : \Phi_2)^{-1}$, and,
 (ii) $\lambda^\downarrow(\Phi_2 : \Phi_1) = \lambda^\uparrow(\Phi_1 : \Phi_2)^{-1}$.

Proof. (i) By proposition 26 $\lambda(\Phi_2 : \Phi_2) = \lambda(\Phi_2, \Phi_1)$. By proposition 24 $\lambda(\Phi_2, \Phi_1) = \lambda(\Phi_1, \Phi_2)^{-1}$, this proves (ii). The map $t \mapsto \frac{1}{t}$ being order reversing, $\lambda^\downarrow(\Phi_2 : \Phi_1) = \lambda^\uparrow(\Phi_1 : \Phi_2)^{-1}$. This proves (iii). \square

E.5 CONCAVITY OF THE LOG-DETERMINANT

Several proofs are possible. We give one leveraging Gaussian integrals and Holder's inequality.

Proposition 28. *The function, $\mathbf{S} \mapsto \log \det [\mathbf{S}]$ is, non-decreasing and concave on the cone of symmetric positive definite matrices.*

Proof. For any $\mathbf{C} \in \mathcal{S}_{>}(d)$, we have, $\int e^{-\frac{1}{2}\langle \mathbf{x}, \mathbf{C} \mathbf{x} \rangle} d\mathbf{x} = \frac{\pi^{\frac{d}{2}}}{\sqrt{\det[\mathbf{C}]}}$. Take $\mathbf{A}, \mathbf{B} \in \mathcal{S}_{>}(d)$, and, $p, q > 0$, such that $\frac{1}{p} + \frac{1}{q} = 1$. We have,

$$\frac{\pi^{\frac{d}{2}}}{\sqrt{\det \left[\frac{1}{p} \mathbf{A} + \frac{1}{q} \mathbf{B} \right]}} = \int e^{-\langle \mathbf{x}, \left(\frac{\mathbf{A}}{p} + \frac{\mathbf{B}}{q} \right) \mathbf{x} \rangle} d\mathbf{x}.$$

By Holder's inequality,

$$\begin{aligned} \int e^{-\langle \mathbf{x}, \frac{\mathbf{A}}{p} \mathbf{x} \rangle} e^{-\langle \mathbf{x}, \frac{\mathbf{B}}{q} \mathbf{x} \rangle} d\mathbf{x} &\leq \left(\int e^{-\langle \mathbf{x}, \mathbf{A} \mathbf{x} \rangle} d\mathbf{x} \right)^{\frac{1}{p}} \left(\int e^{-\langle \mathbf{x}, \mathbf{B} \mathbf{x} \rangle} d\mathbf{x} \right)^{\frac{1}{q}}, \\ &= \left(\frac{\pi^{\frac{d}{2}}}{\sqrt{\det [\mathbf{A}]}} \right)^{\frac{1}{p}} + \left(\frac{\pi^{\frac{d}{2}}}{\sqrt{\det [\mathbf{B}]}} \right)^{\frac{1}{q}}. \end{aligned}$$

Taking the logarithm and simplifying, we conclude,

$$\log \det \left[\frac{1}{p} \mathbf{A} + \frac{1}{q} \mathbf{B} \right] \geq \frac{1}{p} \log \det [\mathbf{A}] + \frac{1}{q} \log \det [\mathbf{B}].$$

\square

2160
2161
2162
2163
2164
2165
2166
2167
2168
2169
2170
2171
2172
2173
2174
2175
2176
2177
2178
2179
2180
2181
2182
2183
2184
2185
2186
2187
2188
2189
2190
2191
2192
2193
2194
2195
2196
2197
2198
2199
2200
2201
2202
2203
2204
2205
2206
2207
2208
2209
2210
2211
2212
2213

E.6 NECESSARY AND SUFFICIENT SOLUTIONS FOR A PSEUDO-INVERSE SOLUTION

Lemma 7 (Penrose 1955). *A necessary and sufficient condition for the equation $\alpha\omega\beta = \gamma$ to have solution is*

$$\alpha\alpha^\dagger\gamma\beta\beta^\dagger = \gamma,$$

in which case the general solution is,

$$\omega = \alpha\alpha^\dagger\gamma\beta\beta^\dagger + \tau - \alpha\alpha^\dagger\alpha\tau\beta\beta^\dagger,$$

E.7 F-DIVERGENCE

Definition 9. Let \mathbb{P} and \mathbb{Q} be two probability measures on a measurable space (Ω, \mathcal{F}) , such that $\mathbb{P} \ll \mathbb{Q}$. For any convex function $f : (0, \infty) \mapsto \mathbb{R}$, strictly convex at 1 and satisfying $f(1) = 0$, the f -divergence between \mathbb{P} and \mathbb{Q} is defined as,

$$D_f(\mathbb{P} \parallel \mathbb{Q}) = \mathbb{E}_{\mathbb{Q}} \left[f \left(\frac{d\mathbb{P}}{d\mathbb{Q}} \right) \right].$$

Theorem 5 (Elementary properties of f -divergences). *Let \mathbb{P} and \mathbb{Q} be two probability measures on the measurable space $(\mathcal{X}, \mathcal{B}(\mathcal{X}))$.*

- (i) $D_f(\mathbb{P} \parallel \mathbb{Q}) \geq 0$ with equality if and only if $\mathbb{P} = \mathbb{Q}$.
- (ii) The function $(\mathbb{P}, \mathbb{Q}) \mapsto D_f(\mathbb{P} \parallel \mathbb{Q})$ is jointly convex.
- (iii) Let $(\mathcal{X}, \mathcal{B}(\mathcal{X}))$ and $(\mathcal{Y}, \mathcal{B}(\mathcal{Y}))$ be two measurable spaces. Let \mathbb{P}_X be a probability measure on $(\mathcal{X}, \mathcal{B}(\mathcal{X}))$. Let $\mathbb{P}_{Y|X} : (\mathcal{X}, \mathcal{B}(\mathcal{X})) \mapsto$ and $\mathbb{Q}_{Y|X}$ be two conditional probability operators (Markov kernels). Define, $\mathbb{P}_Y = \mathbb{E}_{\mathbb{P}_X}[\mathbb{P}_{Y|X}]$ and $\mathbb{Q}_Y = \mathbb{E}_{\mathbb{P}_X}[\mathbb{Q}_{Y|X}]$. Then,

$$D_f(\mathbb{P}_Y \parallel \mathbb{Q}_Y) \leq D_f(\mathbb{P}_{Y|X} \parallel \mathbb{Q}_{Y|X}).$$

Proof. (Polyanskiy & Wu, 2016) □

2214
2215
2216
2217
2218
2219
2220
2221
2222
2223
2224
2225
2226
2227
2228
2229
2230
2231
2232
2233
2234
2235
2236
2237
2238
2239
2240
2241
2242
2243
2244
2245
2246
2247
2248
2249
2250
2251
2252
2253
2254
2255
2256
2257
2258
2259
2260
2261
2262
2263
2264
2265
2266
2267