

# Probabilistic Symmetries and Invariant Neural Networks

**Benjamin Bloem-Reddy**

*Department of Statistics  
University of British Columbia  
Vancouver V6T 1Z4, Canada*

BENBR@STAT.UBC.CA

**Yee Whye Teh**

*Department of Statistics  
University of Oxford  
Oxford OX1 3LB, United Kingdom*

Y.W.TEH@STATS.OX.AC.UK

**Editor:** Ruslan Salakhutdinov

## Abstract

Treating neural network inputs and outputs as random variables, we characterize the structure of neural networks that can be used to model data that are invariant or equivariant under the action of a compact group. Much recent research has been devoted to encoding invariance under symmetry transformations into neural network architectures, in an effort to improve the performance of deep neural networks in data-scarce, non-i.i.d., or unsupervised settings. By considering group invariance from the perspective of probabilistic symmetry, we establish a link between functional and probabilistic symmetry, and obtain generative functional representations of probability distributions that are invariant or equivariant under the action of a compact group. Our representations completely characterize the structure of neural networks that can be used to model such distributions and yield a general program for constructing invariant stochastic or deterministic neural networks. We demonstrate that examples from the recent literature are special cases, and develop the details of the general program for exchangeable sequences and arrays.

**Keywords:** probabilistic symmetry, convolutional neural networks, exchangeability, neural architectures, invariance, equivariance, sufficiency, adequacy, graph neural networks

## 1. Introduction

Neural networks and deep learning methods have found success in a wide variety of applications. Much of the success has been attributed to a confluence of trends, including the increasing availability of data; advances in specialized hardware such as GPUs and TPUs; and open-source software like THEANO (Theano Development Team, 2016), TENSORFLOW (Abadi et al., 2015), and PYTORCH (Paszke et al., 2019), that enable rapid development of neural network models through automatic differentiation and high-level interfaces with specialized hardware. Neural networks have been most successful in settings with massive amounts of i.i.d. labeled training data. In recent years, a concerted research effort has aimed to improve the performance of deep learning systems in data-scarce and semi-supervised or unsupervised problems, and for structured, non-i.i.d. data. In that effort, there has been a renewed focus on novel neural network architectures: attention mechanisms (Vaswani et al.,

2017), memory networks (Sukhbaatar et al., 2015), dilated convolutions (Yu and Koltun, 2016), residual networks (He et al., 2016), and graph neural networks (Scarselli et al., 2009) are a few recent examples from the rapidly expanding literature.

The focus on novel architectures reflects a basic fact of machine learning: in the presence of data scarcity, whether due to small sample size or unobserved variables, or to complicated structure, the model must pick up the slack. Amid the flurry of model-focused innovation, there is a growing need for a framework for encoding modeling assumptions and checking their validity, and for assessing the training stability and generalization potential of an architecture. In short, a principled theory of neural network design is needed.

This paper represents one small step in that direction. It concerns the development of neural network architectures motivated by symmetry considerations, typically described by invariance or equivariance with respect to the action of a group. The most well-known examples of such architectures are convolutional neural networks (CNNs) (LeCun et al., 1989), which ensure invariance of the output  $Y$  under translations of an input image  $X$ . Other examples include neural networks that encode rotational invariance (Cohen et al., 2018) or permutation invariance (Zaheer et al., 2017; Hartford et al., 2018; Herzig et al., 2018; Lee et al., 2019).

Within this growing body of work, symmetry is most often addressed by designing a specific neural network architecture for which the relevant invariance can be verified. A more general approach aims to answer the question:

*For a particular symmetry property, can all invariant neural network architectures be characterized?*

General results have been less common in the literature; important exceptions include characterizations of feed-forward networks (i.e., linear maps composed with pointwise nonlinearities) that are invariant under the action of discrete groups (Shawe-Taylor, 1989; Ravanbakhsh et al., 2017), finite linear groups (Wood and Shawe-Taylor, 1996), or compact groups (Kondor and Trivedi, 2018).

In the probability and statistics literature, there is a long history of probabilistic model specification motivated by symmetry considerations. In particular, if a random variable  $X$  is to be modeled as respecting some symmetry property, formalized as invariance under certain transformations, then a model should only contain invariant distributions  $P_X$ . The relevant theoretical question is:

*For a particular symmetry property of  $X$ , can all invariant distributions  $P_X$  be characterized?*

Work in this area dates at least to the 1930s, and the number of general results reflects the longevity of the field. The most famous example is de Finetti’s theorem (de Finetti, 1930), which is a cornerstone of Bayesian statistics and machine learning. It shows that all infinitely exchangeable (i.e., permutation-invariant) sequences of random variables have distributional representations that are conditionally i.i.d., conditioned on a random probability measure. Other examples include rotational invariance, translation invariance, and a host of others.

In the present work, we approach the question of invariance in neural network architectures from the perspective of probabilistic symmetry. In particular, we seek to understand the symmetry properties of joint probability distributions of  $(X, Y)$  that are necessary and

sufficient for the existence of a *noise-outsourced functional representation*  $Y = f(\eta, X)$ , with generic noise variable  $\eta$ , such that  $f$  obeys the relevant functional symmetries. Our approach sheds light on the core statistical issues involved, and provides a broader view of the questions posed above: from considering classes of deterministic functions to stochastic ones; and from invariant marginal distributions to invariant joint and conditional distributions. Taking a probabilistic approach also leads, in many cases, to simpler proofs of results relative to their deterministic counterparts. Furthermore, stochastic networks are desirable for a number of practical reasons; we discuss these in detail in Section 5.

*Outline.* The remainder of this section provides further background and related work, gives a high-level overview of our main results, and introduces the necessary measure theoretic technicalities and notation. Section 2 defines the relevant functional and probabilistic notions of symmetries; the similarities between the two suggests that there is deeper a mathematical link. Section 3 provides statistical background, reviewing the ideas of sufficiency and adequacy, and of a technical tool known as noise outsourcing. Section 4 makes the precise mathematical link alluded to in Section 2 by establishing functional representations of conditional distributions that are invariant or equivariant under the action of a compact group. Section 5 examines a number of practical considerations and related ideas in the context of the current machine learning literature, and gives a general program for designing invariant neural networks in Section 5.1. Sections 6 and 7 develop the details of the program for exchangeable sequences, arrays, and graphs. Technical details and proofs not given in the main text are in Appendices A to C.

### 1.1. Symmetry in Deep Learning

Interest in neural networks that are invariant to discrete groups acting on the network nodes dates back at least to the text of Minsky and Papert (1988) on single-layer perceptrons (SLPs), who used their results to demonstrate a certain limitation of SLPs. Shawe-Taylor (1989, 1993) extended the theory to multi-layer perceptrons, under the name Symmetry Networks. The main findings of that theory, that invariance is achieved by weight-preserving automorphisms of the neural network, and that the connections between layers must be partitioned into weight-sharing orbits, were rediscovered by Ravanbakhsh et al. (2017), who also proposed novel architectures and new applications.

Wood and Shawe-Taylor (1996) extended the theory to invariance of feed-forward networks under the action of finite linear groups. Some of their results overlap with results for compact groups found in Kondor and Trivedi (2018), including the characterization of equivariance in feed-forward networks in terms of group theoretic convolution.

The most widely applied invariant neural architecture is the CNN for input images. Recently, there has been a surge of interest in generalizing the idea of invariant architectures to other data domains such as sets and graphs, with most work belonging to either of two categories:

- (i) properly defined convolutions (Bruna et al., 2014; Duvenaud et al., 2015; Niepert et al., 2016); or
- (ii) equivariance under the action of groups that lead to weight-tying schemes (Gens and Domingos, 2014; Cohen and Welling, 2016; Ravanbakhsh et al., 2017).

Both of these approaches rely on group theoretic structure in the set of symmetry transformations, and Kondor and Trivedi (2018) used group theory to show that the two approaches are the same (under homogeneity conditions on the group’s action) when the network layers consist of pointwise non-linearities applied to linear maps; Cohen et al. (2019) extended the theory to more general settings. Ravanbakhsh et al. (2017) demonstrated an exception (violating the homogeneity condition) to this correspondence for discrete groups.

Specific instantiations of invariant architectures abound in the literature; they are too numerous to collect here. However, we give a number of examples in Section 5, and in Sections 6 and 7 in the context of sequence- and graph-valued input data that are invariant under permutations.

## 1.2. Symmetry in Probability and Statistics

The study of probabilistic symmetries has a long history. Laplace’s “rule of succession” dates to 1774; it is the conceptual precursor to exchangeability (see Zabell, 2005, for a historical and philosophical account). Other examples include invariance under rotation and stationarity in time (Freedman, 1963); the former has roots in Maxwell’s work in statistical mechanics in 1875 (see, for example, the historical notes in Kallenberg, 2005). The present work relies on the deep connection between sufficiency and symmetry; Diaconis (1988) gives an accessible overview.

In Bayesian statistics, the canonical probabilistic symmetry is exchangeability. A sequence of random variables,  $\mathbf{X}_n = (X_1, \dots, X_n)$ , is exchangeable if its distribution is invariant under all permutations of its elements. If that is true for every  $n \in \mathbb{N}$  in an infinite sequence  $\mathbf{X}_\infty$ , then the sequence is said to be *infinitely exchangeable*. de Finetti’s theorem (de Finetti, 1930; Hewitt and Savage, 1955) shows that infinitely exchangeable distributions have particularly simple structure. Specifically,  $\mathbf{X}_\infty$  is infinitely exchangeable if and only if there exists some random distribution  $Q$  such that the elements of  $\mathbf{X}_\infty$  are conditionally i.i.d. with distribution  $Q$ . Therefore, each infinitely exchangeable distribution  $P$  has an integral decomposition: there is a unique (to  $P$ ) distribution  $\nu$  on the set  $\mathcal{M}_1(\mathcal{X})$  of all probability measures on  $\mathcal{X}$ , such that

$$P(\mathbf{X}_\infty) = \int_{\mathcal{M}_1(\mathcal{X})} \prod_{i=1}^{\infty} Q(X_i) \nu(dQ). \quad (1)$$

The simplicity is useful: by assuming the data are infinitely exchangeable, only models that have a conditionally i.i.d. structure need to be considered. This framework is widely adopted throughout Bayesian statistics and machine learning.

de Finetti’s theorem is a special case of a more general mathematical result, the ergodic decomposition theorem, which puts probabilistic symmetries in correspondence with integral decompositions like (1); see Orbanz and Roy (2015) for an accessible overview. de Finetti’s results inspired a large body of work on other symmetries in the probability literature; Kallenberg (2005) gives a comprehensive treatment and Kallenberg (2017) contains further results. Applications of group symmetries in statistics include equivariant estimation and testing; see, for example, Lehmann and Romano (2005, Ch. 6); Eaton (1989); Wijsman (1990); Giri (1996). The connection between symmetry and statistical sufficiency, which we review in Section 3, was used extensively by Diaconis and Freedman (1984) and

by Lauritzen (1984) in their independent developments of “partial exchangeability” and “extremal families”, respectively; see Diaconis (1988) for an overview. Section 4 makes use of maximal invariants to induce conditional independence; related ideas in a different context were developed by Dawid (1985).

### 1.3. Overview of Main Results

Our main results put functional symmetries in correspondence with probabilistic symmetries. To establish this correspondence, we obtain functional representations of the conditional distribution of an output random variable,  $Y \in \mathcal{Y}$ , given an input random variable,  $X \in \mathcal{X}$ , subject to symmetry constraints. The results provide a general program for constructing deterministic or stochastic functions in terms of statistics that encode a symmetry between input and output, with neural networks as a special case.

*Permutation-invariant output and exchangeable input sequences.* Consider the example of an exchangeable input sequence  $\mathbf{X}_n$ . That is,  $\pi \cdot \mathbf{X}_n \stackrel{d}{=} \mathbf{X}_n$  for all permutations  $\pi \in \mathbb{S}_n$ , where ‘ $\stackrel{d}{=}$ ’ denotes equality in distribution. By definition, the order of the elements has no statistical relevance; only the values contain any information about the distribution of  $\mathbf{X}_n$ . The *empirical measure*<sup>1</sup> (or counting measure),

$$\mathbb{M}_{\mathbf{X}_n}(\cdot) := \sum_{i=1}^n \delta_{X_i}(\cdot),$$

acts to separate relevant from irrelevant statistical information in an exchangeable sequence: it retains the values appearing in  $\mathbf{X}_n$ , but discards their order. In other words, the empirical measure is a *sufficient statistic* for models comprised of exchangeable distributions.

Under certain conditions on the output  $Y$ , the empirical measure also contains all relevant information for predicting  $Y$  from  $\mathbf{X}_n$ , a property known as *adequacy* (see Section 3). In particular, we show in Section 6 that for an exchangeable input  $\mathbf{X}_n$ , the conditional distribution of an output  $Y$  is invariant to permutations of  $\mathbf{X}_n$  if and only if there exists a function  $f$  such that<sup>2</sup>

$$(\mathbf{X}_n, Y) \stackrel{\text{a.s.}}{=} (\mathbf{X}_n, f(\eta, \mathbb{M}_{\mathbf{X}_n})) \quad \text{where } \eta \sim \text{Unif}[0, 1] \quad \text{and } \eta \perp\!\!\!\perp \mathbf{X}_n. \quad (2)$$

This is an example of a noise-outsourced functional representation of samples from a conditional probability distribution. It can be viewed as a general version of the so-called “reparameterization trick” (Kingma and Welling, 2014; Rezende et al., 2014) for random variables taking values in general measurable spaces (not just  $\mathbb{R}$ ). The noise variable  $\eta$  acts as a generic source of randomness that is “outsourced”, a term borrowed from Austin (2015). The relevant property of  $\eta$  is its independence from  $X$ , and the uniform distribution is not special in this regard; it could be replaced by, for example, a standard normal random variable and the result would still hold, albeit with a different  $f$ . For modeling purposes, the outer function  $f$  may contain “pooling functions” that compress  $\mathbb{M}_{\mathbf{X}_n}$  further, such as sum or max. The only restriction imposed by being a function of  $\mathbb{M}_{\mathbf{X}_n}$  is that the order of

1. The Dirac delta function  $\delta_x(B) = 1$  if  $x \in B$  and is 0 otherwise, for any measurable set  $B$ .

2. ‘ $\stackrel{\text{a.s.}}{=}$ ’ denotes equal almost surely.

elements in  $\mathbf{X}_n$  has been discarded:  $f$  must be invariant to permutations of  $\mathbf{X}_n$ . We give further examples in Section 6.

The representation (2) is a characterization of permutation-invariant stochastic functions, and generalizes a deterministic version that Zaheer et al. (2017) obtained under more restrictive assumptions; Murphy et al. (2019) and Wagstaff et al. (2019), among others, provide extensions and further theoretical study of deterministic permutation-invariant functions.

*Permutation-equivariant output and exchangeable input sequences.* For the purposes of constructing deep neural networks, invariance often imposes stronger constraints than are desirable for hidden layers. Instead, the output should transform predictably under transformations of the input, in a manner that preserves the structure of the input; this is a property known as *equivariance*. For simplicity, consider an output,  $\mathbf{Y}_n = (Y_1, \dots, Y_n)$ , of the same size as the input  $\mathbf{X}_n$ . For exchangeable input  $\mathbf{X}_n$ , the conditional distribution of  $\mathbf{Y}_n$  is equivariant to permutations of  $\mathbf{X}_n$  if and only if there exists a function  $f$  such that<sup>3</sup>

$$(\mathbf{X}_n, \mathbf{Y}_n) \stackrel{\text{a.s.}}{=} (\mathbf{X}_n, (f(\eta_i, X_i, \mathbb{M}_{\mathbf{X}_n}))_{1 \leq i \leq n}) \quad \text{where } \eta_i \stackrel{\text{iid}}{\sim} \text{Unif}[0, 1] \quad \text{and } (\eta_i)_{1 \leq i \leq n} \perp\!\!\!\perp \mathbf{X}_n .$$

Special (deterministic) cases of  $f$  have appeared in the literature (Zaheer et al., 2017; Lee et al., 2019), with equivariance demonstrated on a case-by-case basis; to the best of our knowledge, no version of this general result has previously appeared. We give further examples in Section 6.

*Invariance and equivariance under compact groups.* The representation (2) explicitly separates the conditional distribution of  $Y$  given  $\mathbf{X}_n$  into relevant structure (the empirical measure  $\mathbb{M}_{\mathbf{X}_n}$ ) and independent random noise. Such separation is a recurring theme throughout the present work. For more general groups of transformations than permutations, the relevant structure is captured by a *maximal invariant* statistic  $M(X)$ , which can be used to partition the input space into equivalence classes under the action of a group  $\mathcal{G}$ . We show in Section 4 that for an input  $X$  satisfying  $g \cdot X \stackrel{\text{d}}{=} X$  for all  $g \in \mathcal{G}$ , the conditional distribution of an output  $Y$  is invariant to the action of  $\mathcal{G}$  on  $X$  if and only if there exists some function  $f$  such that

$$(X, Y) \stackrel{\text{a.s.}}{=} (X, f(\eta, M(X))) \quad \text{where } \eta \sim \text{Unif}[0, 1] \quad \text{and } \eta \perp\!\!\!\perp X . \quad (3)$$

Because  $M$  is  $\mathcal{G}$ -invariant and  $f$  depends on  $X$  only through  $M$ ,  $f$  is also  $\mathcal{G}$ -invariant. The empirical measure is an example of a maximal invariant for  $\mathbb{S}_n$  acting on  $\mathcal{X}^n$ , leading to (2) as a special case of (3).

We also obtain a functional representation of all  $\mathcal{G}$ -equivariant conditional distributions: for input  $X$  with  $\mathcal{G}$ -invariant distribution, the conditional distribution of  $Y$  given  $X$  is  $\mathcal{G}$ -equivariant if and only if

$$(X, Y) \stackrel{\text{a.s.}}{=} (X, f(\eta, X)) \quad \text{where } \eta \sim \text{Unif}[0, 1] \quad \text{and } \eta \perp\!\!\!\perp X ,$$

for  $f$  satisfying

$$g \cdot Y = g \cdot f(\eta, X) = f(\eta, g \cdot X) , \quad \text{a.s., for all } g \in \mathcal{G} .$$

---

3. The representation also requires that the elements of  $\mathbf{Y}_n$  are conditionally independent given  $\mathbf{X}_n$ , which is trivially satisfied by most neural network architectures; see Section 6.2.

Specifically, distributional equivariance is equivalent to for  $(X, Y) \stackrel{d}{=} (g \cdot X, g \cdot Y)$ ,  $g \in \mathcal{G}$ . For  $(X, Y)$  satisfying that condition, one may construct an equivariant function from a *representative equivariant*  $\tau(X)$ , which is an element of  $\mathcal{G}$  that maps  $X$  to a representative element of its equivalence class. See Section 4.2 for details.

*Exchangeable arrays and graphs.* In addition to applying the general theory to exchangeable sequences (Section 6), we develop the details of functional representations of exchangeable arrays and graphs in Section 7. Consider an input array  $\mathbf{X}_{n_2}$ , modeled as invariant under separate permutations of its rows and columns. Such *separately exchangeable* arrays arise in, for example, collaborative filtering problems like recommender systems. The analogue of the empirical measure in (2) is any *canonical form*  $\mathbf{C}_{\mathbf{X}_{n_2}} := \text{CANON}(\mathbf{X}_{n_2})$ . The conditional distribution of an output  $Y$  is invariant under permutations of the rows and columns of a separately exchangeable input  $\mathbf{X}_{n_2}$  if and only if there exists a function  $f$  such that

$$(\mathbf{X}_{n_2}, Y) \stackrel{\text{a.s.}}{=} (\mathbf{X}_{n_2}, f(\eta, \mathbf{C}_{\mathbf{X}_{n_2}})) \quad \text{where } \eta \sim \text{Unif}[0, 1] \quad \text{and } \eta \perp\!\!\!\perp \mathbf{X}_{n_2} .$$

Analogously, the conditional distribution of an output array  $\mathbf{Y}_{n_2}$  is equivariant under permutations of the rows and columns of  $\mathbf{X}_{n_2}$  if and only if there exists a function  $f$  such that

$$(\mathbf{X}_{n_2}, \mathbf{Y}_{n_2}) \stackrel{\text{a.s.}}{=} (\mathbf{X}_{n_2}, (f(\eta_{i,j}, X_{i,j}, \mathbf{C}_{\mathbf{X}_{n_2}}^{(i,j)}))_{i \leq n_1, j \leq n_2}),$$

where  $\mathbf{C}_{\mathbf{X}_{n_2}}^{(i,j)}$  is an augmented version of  $\mathbf{C}_{\mathbf{X}_{n_2}}$ , with the  $i$ th row and  $j$ th column of  $\mathbf{X}_{n_2}$  broadcast to the entire array. Deterministic special cases of these representations have appeared in, for example Hartford et al. (2018); Herzig et al. (2018). See Section 7 for details and more examples, where analogous results are developed for graphs (i.e., symmetric arrays that are exchangeable under the same permutation applied to the rows and columns) and for arrays with features associated with each row and column. Results for  $d$ -dimensional arrays (i.e., tensors) are given in Appendix C.

*Contributions.* The probabilistic techniques employed here are generally not new (most are fairly standard), though to our knowledge they have not been applied to the types of problems considered here. Furthermore, the invariant deep learning literature has been in apparent disconnect from important ideas concerning probabilistic symmetry, particularly as related to statistical sufficiency (see Section 3.1) and foundational ideas in Bayesian statistics; and vice versa. This allows for exchange of ideas between the areas. In addition to certain practical advantages of taking a probabilistic approach (which we detail in Section 5.3), we obtain *the most general possible representation results* for stochastic or deterministic functions that are invariant or equivariant with respect to a compact group. Our representation results are *exact*, as opposed to recent results on universal approximation (see Section 5.4). To the best of our knowledge, no results at that level of generality exist in the deep learning literature.

#### 1.4. Technicalities and Notation

In order to keep the presentation as clear as possible, we aim to minimize measure theoretic technicalities. Throughout, it is assumed that there is a background probability space

$(\Omega, \mathcal{A}, \mathbb{P})$  that is rich enough to support all required random variables. All random variables are assumed to take values in standard Borel spaces (spaces that are Borel isomorphic to a Borel subset of the unit interval (see, e.g., Kallenberg, 2002)). For example,  $X$  is a  $\mathcal{X}$ -valued random variable in  $(\mathcal{X}, \mathcal{B}_{\mathcal{X}})$ , where  $\mathcal{B}_{\mathcal{X}}$  is the Borel  $\sigma$ -algebra of  $\mathcal{X}$ . Alternatively, we may say that  $X$  is a random element of  $\mathcal{X}$ . For notational convenience, for a  $\mathcal{Y}$ -valued random variable  $Y$ , we write  $P(Y \in \bullet \mid X)$  as shorthand for, “for all sets  $A \in \mathcal{B}_{\mathcal{Y}}$ ,  $P(Y \in A \mid \sigma(X))$ ”, where  $\sigma(X)$  is the  $\sigma$ -algebra generated by  $X$ . We use  $\mathcal{M}(\mathcal{X})$  to denote the set of measures on  $\mathcal{X}$ , and  $\mathcal{M}_1(\mathcal{X})$  to denote the set of probability measures on  $\mathcal{X}$ . Many of our results pertain to conditional independence relationships;  $Y \perp\!\!\!\perp_Z X$  means that  $Y$  and  $X$  are conditionally independent, given  $\sigma(Z)$ . Finally,  $\stackrel{d}{=}$  denotes equality in distribution, and  $\stackrel{\text{a.s.}}{=}$  denotes almost sure equality.

We consider symmetries induced by the action of a group. A group is a set,  $\mathcal{G}$ , and a binary composition operator  $\cdot$  that together must satisfy four properties: for each  $g, g' \in \mathcal{G}$  the composition  $g \cdot g' \in \mathcal{G}$  is in the group; the group operation is associative, that is,  $g \cdot (g' \cdot g'') = (g \cdot g') \cdot g''$  for all  $g, g', g'' \in \mathcal{G}$ ; there is an identity element  $e \in \mathcal{G}$  such that  $e \cdot g = g \cdot e = g$  for all  $g \in \mathcal{G}$ ; and for each  $g \in \mathcal{G}$  there is an inverse  $g^{-1} \in \mathcal{G}$  such that  $g^{-1} \cdot g = g \cdot g^{-1} = e$ . See, e.g., Rotman (1995). Let  $\Phi_{\mathcal{X}} : \mathcal{G} \times \mathcal{X} \rightarrow \mathcal{X}$  be the left-action of  $\mathcal{G}$  on the input space  $\mathcal{X}$ , such that  $\Phi_{\mathcal{X}}(e, x) = x$  is the identity mapping for all  $x \in \mathcal{X}$ , and  $\Phi_{\mathcal{X}}(g, \Phi_{\mathcal{X}}(g', x)) = \Phi_{\mathcal{X}}(g \cdot g', x)$  for  $g, g' \in \mathcal{G}$ ,  $x \in \mathcal{X}$ . For convenience, we write  $g \cdot x = \Phi_{\mathcal{X}}(g, x)$ . Similarly, let  $\Phi_{\mathcal{Y}}$  be the action of  $\mathcal{G}$  on the output space  $\mathcal{Y}$ , and  $g \cdot y = \Phi_{\mathcal{Y}}(g, y)$  for  $g \in \mathcal{G}$ ,  $y \in \mathcal{Y}$ . In this paper, we always assume  $\mathcal{G}$  to be compact. A group  $\mathcal{G}$ , along with a  $\sigma$ -algebra  $\sigma(\mathcal{G})$ , is said to be measurable if the group operations of inversion  $g \mapsto g^{-1}$  and composition  $(g, g') \mapsto g \cdot g'$  are  $\sigma(\mathcal{G})$ -measurable.  $\mathcal{G}$  acts measurably on  $\mathcal{X}$  if  $\Phi_{\mathcal{X}}$  is a measurable function  $\sigma(\mathcal{G}) \otimes \mathcal{B}_{\mathcal{X}} \rightarrow \mathcal{B}_{\mathcal{X}}$  (Kallenberg, 2017).

## 2. Functional and Probabilistic Symmetries

We consider the relationship between two random variables  $X \in \mathcal{X}$  and  $Y \in \mathcal{Y}$ , with  $Y$  being a predicted output based on input  $X$ . For example, in image classification,  $X$  might be an image and  $Y$  a class label; in sequence prediction,  $X = (X_i)_{i=1}^n$  might be a sequence and  $Y$  the next element,  $X_{n+1}$ , to be predicted; and in a variational autoencoder,  $X$  might be an input vector and  $Y$  the corresponding latent variable whose posterior is to be inferred using the autoencoder. Throughout, we denote the joint distribution of both variables as  $P_{X,Y}$ , and the conditional distribution  $P_{Y|X}$  is the primary object of interest.

Two basic notions of symmetry are considered, as are the connections between them. The first notion, defined in Section 2.1, is functional, and is most relevant when  $Y$  is a deterministic function of  $X$ , say  $Y = f(X)$ ; the symmetry properties pertain to the function  $f$ . The second notion, defined in Section 2.2, is probabilistic, and pertains to the conditional distribution of  $Y$  given  $X$ .

### 2.1. Functional Symmetry in Neural Networks

In many machine learning settings, a prediction  $y$  based on an input  $x$  is modeled as a deterministic function,  $y = f(x)$ , where  $f$  belongs to some function class  $\mathcal{F} = \{f; f : \mathcal{X} \rightarrow \mathcal{Y}\}$ , often satisfying some further conditions. For example,  $f$  might belong to a Reproducing Kernel Hilbert Space, or to a subspace of all strongly convex functions. Alternatively,  $f$



may be a neural network parameterized by weights and biases collected into a parameter vector  $\theta$ , in which case the function class  $\mathcal{F}$  corresponds to the chosen network architecture, and a particular  $f$  corresponds to a particular set of values for  $\theta$ . We are concerned with implications on the choice of the network architecture (equivalently, the function class  $\mathcal{F}$ ) due to symmetry properties we impose on the input-output relationship. In this deterministic setting, the conditional distribution  $P_{Y|X}$  is simply a point mass at  $f(x)$ .

Two properties, invariance and equivariance, formalize the relevant symmetries. A function  $f : \mathcal{X} \rightarrow \mathcal{Y}$  is *invariant under  $\mathcal{G}$* , or  *$\mathcal{G}$ -invariant*, if the output is unchanged by transformations of the input induced by the group:

$$f(g \cdot x) = f(x) \quad \text{for all } g \in \mathcal{G}, x \in \mathcal{X} . \quad (4)$$

Alternatively, a function  $f : \mathcal{X} \rightarrow \mathcal{Y}$  is *equivariant under  $\mathcal{G}$* , or  *$\mathcal{G}$ -equivariant*, if

$$f(g \cdot x) = g \cdot f(x) \quad \text{for all } g \in \mathcal{G}, x \in \mathcal{X} . \quad (5)$$

The action of  $\mathcal{G}$  commutes with the application of an equivariant  $f$ ; transforming the input is the same as transforming the output. Note that the action of  $\mathcal{G}$  on  $\mathcal{X}$  and  $\mathcal{Y}$  may be different. In particular, invariance is a special case of equivariance, whereby the group action in the output space is trivial:  $g \cdot y = y$  for each  $g \in \mathcal{G}$  and  $y \in \mathcal{Y}$ . Invariance imposes stronger restrictions on the functions satisfying it, as compared to equivariance.

These properties and their implications for network architectures are illustrated with examples from the literature. For notational convenience, let  $[n] = \{1, \dots, n\}$  and  $\mathbf{X}_n = (X_1, \dots, X_n) \in \mathcal{X}^n$ . Finally, denote the finite symmetric group of a set of  $n$  elements (i.e., the set of all permutations of  $[n]$ ) by  $\mathbb{S}_n$ .

**Example 1 (Deep Sets:  $\mathbb{S}_n$ -invariant functions of sequences)** *Zaheer et al. (2017) considered a model  $Y = f(\mathbf{X}_n)$ , where the input  $\mathbf{X}_n$  was treated as a set, i.e., the order among its elements did not matter. Those authors required that the output of  $f$  be unchanged under all permutations of the elements of  $\mathbf{X}_n$ , i.e., that  $f$  is  $\mathbb{S}_n$ -invariant. They found that  $f$  is  $\mathbb{S}_n$ -invariant if and only if it can be represented as  $f(\mathbf{X}_n) = \tilde{f}(\sum_{i=1}^n \phi(X_i))$  for some functions  $\tilde{f}$  and  $\phi$ . Clearly, permutations of the elements of  $\mathbf{X}_n$  leave such a function invariant:  $f(\mathbf{X}_n) = f(\pi \cdot \mathbf{X}_n)$  for  $\pi \in \mathbb{S}_n$ . The fact that all  $\mathbb{S}_n$ -invariant functions can be expressed in such a form was proved by Zaheer et al. (2017) in two different settings: (i) for sets of arbitrary size when  $\mathcal{X}$  is countable; and (ii) for sets of fixed size when  $\mathcal{X}$  is uncountable. In both proofs, the existence of such a form is shown by constructing a  $\phi$  that uniquely encodes the elements of  $\mathcal{X}$ . Essentially,  $\phi$  is a generalization of a one-hot encoding, which gets sum-pooled and passed through  $\tilde{f}$ . The authors call neural architectures satisfying such structure Deep Sets. See Figure 1, left panel, for an example diagram. In Section 6.1, we give a short and intuitive proof using basic properties of exchangeable sequences.*

**Example 2 ( $\mathbb{S}_n$ -equivariant neural network layers)** *Let  $\mathbf{X}_n$  and  $\mathbf{Y}_n$  represent adjacent layers of a standard feed-forward neural network, such that the nodes are indexed by  $[n]$ , with  $X_i \in \mathbb{R}$  the  $i$ th node in layer  $\mathbf{X}_n$ , and similarly for  $Y_i$ . In a feed-forward layer,  $\mathbf{Y}_n = \sigma(\theta \mathbf{X}_n)$  where  $\sigma$  is an element-wise nonlinearity and  $\theta$  is a weight matrix (we ignore biases for simplicity). Shawe-Taylor (1989); Wood and Shawe-Taylor (1996); Zaheer et al.*

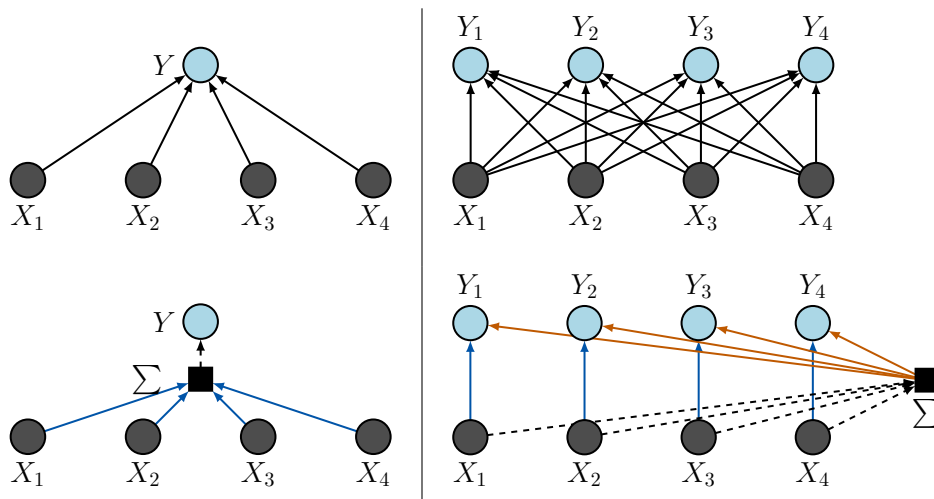


Figure 1: Computation diagrams for the neural networks in Examples 1 and 2. Black solid arrows indicate general, uncoupled weights; black dashed arrows indicate a fixed function, i.e., an activation function; colored arrows indicate shared weights between arrows of the same color. *Left panel:* A general output layer with a different weight from each  $X_i$  to  $Y$  (top), and a simple weight-sharing  $S_n$ -invariant architecture (bottom). *Right panel:* A fully connected MLP with  $n^2$  weights (top), and the  $S_n$ -equivariant architecture corresponding to (6), with two weights (bottom).

(2017) showed that the only weight matrices that lead to  $S_n$ -equivariant layers are the sum of a diagonal matrix,  $\theta_0 \mathbb{I}_n$ ,  $\theta_0 \in \mathbb{R}$ , and a constant one,  $\theta_1 \mathbf{1}_n \mathbf{1}_n^T$ ,  $\theta_1 \in \mathbb{R}$ , such that

$$[\theta \mathbf{X}_n]_i = \theta_0 X_i + \theta_1 \sum_{j=1}^n X_j. \quad (6)$$

Figure 1, right panel, shows the weight-sharing patterns of the connections.

These examples demonstrate that equivariance is less restrictive than invariance, and thus allows for more expressive parameterizations; in Example 2, invariance would require that  $\theta_0 = 0$ . At the same time, equivariance appears to be strong enough to greatly reduce the dimension of the parameter space: generic fully connected layers contain  $n^2$  weights, compared to the two used by the  $S_n$ -equivariant architecture.

At a high level, equivariance in feed-forward neural networks ensures that transformations of the input lead to predictable, symmetry-preserving transformations of higher layers, which allows the network to exploit the symmetry in all layers through weight-sharing (Cohen and Welling, 2016), and to summarize features at multiple scales through pooling (Kondor and Trivedi, 2018). In addition to being theoretically interesting, group invariance often indicates simplified architectures through parameter-sharing (e.g., Ravanbakhsh et al., 2017; Cohen and Welling, 2017). These in turn lead to simpler, more stable training and

may lead to better generalization (Shawe-Taylor, 1991). Related ideas pertaining to the benefits of invariance in the context of data augmentation and feature averaging are found in Chen et al. (2019); Lyle et al. (2020). (See also the discussion in Section 8.)

## 2.2. Symmetry in Conditional Probability Distributions

An alternative approach to the deterministic models of Section 2.1 is to model the relationship between input  $X$  and output  $Y$  as stochastic, either by directly parameterizing the conditional distribution  $P_{Y|X}$ , or by defining a procedure for generating samples from  $P_{Y|X}$ . For example, the encoder network of a variational autoencoder computes an approximate posterior distribution over the latent variable  $Y$  given observation  $X$ ; in classification, the use of a soft-max output layer is interpreted as a network which predicts a distribution over labels; in implicit models like Generative Adversarial Networks (Goodfellow et al., 2014) or simulation-based models without a likelihood (e.g., Gutmann and Corander, 2016), and in many probabilistic programming languages (e.g., van de Meent et al., 2018),  $P_{Y|X}$  is not explicitly represented, but samples are generated and used to evaluate the quality of the model.

The relevant properties that encode symmetry in such settings are therefore probabilistic. One way to define symmetry properties for conditional distributions is by adding noise to invariant or equivariant functions. For example, if  $f$  is  $\mathcal{G}$ -invariant and  $\eta$  is a standard normal random variable independent of  $X$ , then  $Y = f(X) + \eta$  corresponds to a conditional distribution  $P_{Y|X}$  that is  $\mathcal{G}$ -invariant. This type of construction, with  $Y = f(\eta, X)$  satisfying invariance (equivariance) in its second argument, will lead to invariant (equivariant) conditional distributions. While intuitive and constructive, it does not characterize what probabilistic properties must be satisfied by random variables  $X$  and  $Y$  in order for such representations to exist. Furthermore, it leaves open the question of whether there are other approaches that may be used. To avoid these ambiguities, we define notions of symmetries for probability models directly in terms of the distributions.

The discussion on exchangeability in Section 1 pertains only to invariance of the marginal distribution  $P_X$  under  $\mathbb{S}_n$ . Suitable notions of probabilistic symmetry under more general groups are needed for the conditional distribution of  $Y$  given  $X$ . Let  $\mathcal{G}$  be a group acting measurably on  $\mathcal{X}$  and on  $\mathcal{Y}$ . The conditional distribution  $P_{Y|X}$  of  $Y$  given  $X$  is  $\mathcal{G}$ -invariant if  $Y|X \stackrel{\text{d}}{=} Y|g \cdot X$ . More precisely, for all  $A \in \mathcal{B}_Y$  and  $B \in \mathcal{B}_X$ ,

$$P_{Y|X}(Y \in A \mid X \in B) = P_{Y|X}(Y \in A \mid g \cdot X \in B) \quad \text{for all } g \in \mathcal{G}. \quad (7)$$

On the other hand,  $P_{Y|X}$  is  $\mathcal{G}$ -equivariant if  $Y|X \stackrel{\text{d}}{=} g \cdot Y|g \cdot X$  for all  $g \in \mathcal{G}$ . That is, transforming  $X$  by  $g$  leads to the same conditional distribution of  $Y$  except that  $Y$  is also transformed by  $g$ . More precisely, for all  $A \in \mathcal{B}_Y$  and  $B \in \mathcal{B}_X$ ,

$$P_{Y|X}(Y \in A \mid X \in B) = P_{Y|X}(g \cdot Y \in A \mid g \cdot X \in B) \quad \text{for all } g \in \mathcal{G}. \quad (8)$$

Typically, conditional invariance or equivariance is desired because of symmetries in the marginal distribution  $P_X$ . In Example 1, it was assumed that the ordering among the elements of the input sequence is unimportant, and therefore it is reasonable to assume that the marginal distribution of the input sequence is exchangeable. In general, we assume

throughout the present work that  $X$  is marginally  $\mathcal{G}$ -invariant:

$$P_X(X \in B) = P_X(g \cdot X \in B) \quad \text{for all } g \in \mathcal{G}, B \in \mathcal{B}_X. \quad (9)$$

Clearly,  $\mathcal{G}$ -invariance of  $P_X$  and of  $P_{Y|X}$  will result in  $P_{X,Y}(X, Y) = P_{X,Y}(g \cdot X, Y)$ ; similarly, if  $P_X$  is  $\mathcal{G}$ -invariant and  $P_{Y|X}$  is  $\mathcal{G}$ -equivariant, then  $P_{X,Y}(X, Y) = P_{X,Y}(g \cdot X, g \cdot Y)$ . The converse is also true for sufficiently nice groups and spaces (such that conditioning is well-defined). Therefore, we may work with the joint distribution of  $X$  and  $Y$ , which is often more convenient than working with the marginal and conditional distributions. These ideas are summarized in the following proposition, which is a special case of more general results on invariant measures found in, for example, Kallenberg (2017, Ch. 7).

**Proposition 1** *For a group  $\mathcal{G}$  acting measurably on Borel spaces  $\mathcal{X}$  and  $\mathcal{Y}$ , if  $P_X$  is marginally  $\mathcal{G}$ -invariant then*

- (i)  $P_{Y|X}$  is conditionally  $\mathcal{G}$ -invariant if and only if  $(X, Y) \stackrel{d}{=} (g \cdot X, Y)$  for all  $g \in \mathcal{G}$ .
- (ii)  $P_{Y|X}$  is conditionally  $\mathcal{G}$ -equivariant if and only if  $(X, Y) \stackrel{d}{=} g \cdot (X, Y) := (g \cdot X, g \cdot Y)$  for all  $g \in \mathcal{G}$ , i.e.,  $X$  and  $Y$  are jointly  $\mathcal{G}$ -invariant.

There is a simple algebra of compositions of equivariant and invariant functions. Specifically, compositions of equivariant functions are equivariant (equivariance is transitive under function composition), and composing an equivariant function with an invariant one yields an invariant function (Cohen and Welling, 2016; Kondor and Trivedi, 2018). Such compositions generate a grammar with which to construct elaborate functions with the desired symmetry properties. Examples abound in the deep learning literature. These compositional properties carry over to the probabilistic case as well.

**Proposition 2** *Let  $X, Y, Z$  be random variables such that  $X \perp\!\!\!\perp_Y Z$ . Suppose  $X$  is marginally  $\mathcal{G}$ -invariant.*

- (i) *If  $P_{Y|X}$  is conditionally  $\mathcal{G}$ -equivariant, and  $P_{Z|Y}$  is conditionally  $\mathcal{G}$ -equivariant, then  $P_{Z|X}$  is also conditionally  $\mathcal{G}$ -equivariant.*
- (ii) *If  $P_{Y|X}$  is conditionally  $\mathcal{G}$ -equivariant, and  $P_{Z|Y}$  is conditionally  $\mathcal{G}$ -invariant, then  $P_{Z|X}$  is conditionally  $\mathcal{G}$ -invariant.*

**Proof** (i) Theorem 1 shows that  $X, Y$  are jointly  $\mathcal{G}$ -invariant, implying that  $Y$  is marginally  $\mathcal{G}$ -invariant. Theorem 1 applied to  $Y$  and  $Z$  shows that  $Y, Z$  are jointly  $\mathcal{G}$ -invariant as well. Joint invariance of  $X, Y$ , along with  $X \perp\!\!\!\perp_Y Z$ , implies that  $X, Y, Z$  are jointly  $\mathcal{G}$ -invariant. After marginalizing out  $Y$ , the joint distribution of  $(X, Z)$  is  $\mathcal{G}$ -invariant. Theorem 1 then shows that  $P_{Z|X}$  is conditionally  $\mathcal{G}$ -equivariant.

(ii) A similar argument as above shows that  $(g \cdot X, g \cdot Y, Z) \stackrel{d}{=} (X, Y, Z)$  for each  $g \in \mathcal{G}$ . Marginalizing out  $Y$ , we see that  $(g \cdot X, Z) \stackrel{d}{=} (X, Z)$ , so that  $P_{Z|X}$  is conditionally  $\mathcal{G}$ -invariant. ■

### 3. Sufficiency, Adequacy, and Noise Outsourcing

Functional and probabilistic notions of symmetries represent two different approaches to achieving the same goal: a principled framework for constructing models from symmetry considerations. Despite their apparent similarities, the precise mathematical link is not immediately clear. The connection consists of two components: *noise outsourcing*, a generic technical tool, allows us to move between a conditional probability distribution  $P_{Y|X}$  and a representation of  $Y$  as a function of  $X$  and random noise; *statistical sufficiency and adequacy* allow us to restrict the noise-outsourced representations to functions of certain statistics of  $X$ . We review these ideas in this section.

#### 3.1. Sufficiency and Adequacy

Sufficiency is based on the idea that a statistic may contain all information that is needed for an inferential procedure; for example, to completely describe the distribution of a sample or for parameter inference. Adequacy extends that idea to prediction. The ideas go hand-in-hand with notions of symmetry: while invariance describes information that is irrelevant, sufficiency and adequacy describe the information that is relevant.

Sufficiency and adequacy are defined with respect to a probability *model*: a family of distributions indexed by some parameter  $\theta \in \Theta$ . Throughout, we consider a model for the joint distribution over  $X$  and  $Y$ ,  $\mathcal{P}_{X,Y} = \{P_{X,Y;\theta} : \theta \in \Theta\}$ , from which there is an induced marginal model  $\mathcal{P}_X = \{P_{X;\theta} : \theta \in \Theta\}$  and a conditional model  $\mathcal{P}_{Y|X} = \{P_{Y|X;\theta} : \theta \in \Theta\}$ . For convenience, and because the parameter does not play a meaningful role in subsequent developments, we suppress the notational dependence on  $\theta$  when there is no chance of confusion.

Sufficiency originates in the work of Fisher (1922); it formalizes the notion that a statistic might be used in place of the data for any statistical procedure. It has been generalized in a variety of different directions, including predictive sufficiency (Bahadur, 1954; Lauritzen, 1974a; Fortini et al., 2000) and adequacy (Skibinsky, 1967). There are a number of ways to formalize sufficiency, which are equivalent under certain regularity conditions; see Schervish (1995). The definition that is most convenient here is due to Halmos and Savage (1949): there is a *single* Markov kernel that gives the *same* conditional distribution of  $X$  conditioned on  $S(X) = s$  for *every* distribution  $P_X \in \mathcal{P}_X$ .

**Definition 3** *Let  $\mathcal{S}$  be a Borel space and  $S : \mathcal{X} \rightarrow \mathcal{S}$  a measurable map.  $S$  is a sufficient statistic for  $\mathcal{P}_X$  if there is a Markov kernel  $q : \mathcal{B}_X \times \mathcal{S} \rightarrow \mathbb{R}_+$  such that for all  $P_X \in \mathcal{P}_X$  and  $s \in \mathcal{S}$ , we have  $P_X(\cdot | S(X) = s) = q(\cdot, s)$ .*

A canonical example is that of a sequence of  $n$  i.i.d. coin tosses. If the probability model is the family of Bernoulli distributions with probability of heads equal to  $p$ , then the number of heads  $N_h$  is sufficient: conditioned on  $N_h = n_h$ , the distribution of the data is uniform on all sequences with  $n_h$  heads. Equivalently, the number of heads is also sufficient in estimating  $p$ . For example, the maximum likelihood estimator is  $N_h/n$ .

Section 6 pertains to the more nuanced example of finite exchangeable sequences  $\mathbf{X}_n \in \mathcal{X}^n$ . A distribution  $P_X$  on  $\mathcal{X}^n$  is *finitely exchangeable* if for all sets  $A_1, \dots, A_n \in \mathcal{B}_X$ ,

$$P(X_1 \in A_1, \dots, X_n \in A_n) = P(X_{\pi(1)} \in A_1, \dots, X_{\pi(n)} \in A_n), \quad \text{for all } \pi \in \mathbb{S}_n. \quad (10)$$

Denote by  $\mathcal{P}_{\mathbf{X}_n}^{\mathbb{S}_n}$  the family of all exchangeable distributions on  $\mathcal{X}^n$ . The de Finetti conditional i.i.d. representation (1) may fail for a finitely exchangeable sequence (see Diaconis, 1977; Diaconis and Freedman, 1980a, for examples). However, the empirical measure plays a central role in both the finitely and infinitely exchangeable cases. The empirical measure (or counting measure) of a sequence  $\mathbf{X}_n \in \mathcal{X}^n$  is defined as

$$\mathbb{M}_{\mathbf{X}_n}(\bullet) = \sum_{i=1}^n \delta_{X_i}(\bullet), \quad (11)$$

where  $\delta_{X_i}$  denotes an atom of unit mass at  $X_i$ . The empirical measure discards the information about the order of the elements of  $\mathbf{X}_n$ , but retains all other information. A standard fact (see Section 6) is that a distribution  $P_{\mathbf{X}_n}$  on  $\mathcal{X}^n$  is exchangeable if and only if the conditional distribution  $P_{\mathbf{X}_n}(\mathbf{X}_n \mid \mathbb{M}_{\mathbf{X}_n} = m)$  is the uniform distribution on all sequences that can be obtained by applying a permutation to  $\mathbf{X}_n$ . This is true for all distributions in  $\mathcal{P}_{\mathbf{X}_n}^{\mathbb{S}_n}$ ; therefore  $\mathbb{M}_{\mathbf{X}_n}$  is a sufficient statistic for  $\mathcal{P}_{\mathbf{X}_n}^{\mathbb{S}_n}$  according to Theorem 3, and we may conclude for any probability model  $P_{\mathbf{X}_n}$ :

*$\mathbb{S}_n$ -invariance of all  $P_{\mathbf{X}_n} \in \mathcal{P}_{\mathbf{X}_n}$  is equivalent to the sufficiency of  $\mathbb{M}_{\mathbf{X}_n}$ .*

In this case invariance and sufficiency clearly are two sides of the same coin: a sufficient statistic captures all information that is relevant to a model for  $\mathbf{X}_n$ ; invariance discards the irrelevant information. Section 6 explores this in further detail.

We note that in this work, there is a clear correspondence between group invariance and a sufficient statistic (see Section 4.3 for details), as in the example of exchangeable sequences. In other situations, there is a sufficient statistic but the set of symmetry transformations may not correspond to a group. See Freedman (1962, 1963); Diaconis and Freedman (1987) for examples.

*Adequacy.* The counterpart of sufficiency for modeling the conditional distribution of  $Y$  given  $X$  is *adequacy* (Skibinsky, 1967; Speed, 1978).<sup>4</sup> The following definition adapts one given by Lauritzen (1974b), which is more intuitive than the measure theoretic definition introduced by Skibinsky (1967).

**Definition 4** *Let  $\mathcal{S}$  be a Borel space, and let  $S : \mathcal{X} \rightarrow \mathcal{S}$  be a measurable map. Then  $S$  is an adequate statistic of  $X$  for  $Y$  with respect to  $\mathcal{P}_{X,Y}$  if*

- (i)  *$S$  is sufficient for  $\mathcal{P}_X$ ; and*
- (ii) *for all  $x \in \mathcal{X}$  and  $P_{X,Y} \in \mathcal{P}_{X,Y}$ ,*

$$P_{X,Y}(Y \in \bullet \mid X = x) = P_{X,Y}(Y \in \bullet \mid S = S(x)). \quad (12)$$

Equation (12) amounts to the conditional independence of  $Y$  and  $X$ , given  $S(X)$ . To see this, note that because  $S(X)$  is a measurable function of  $X$ ,

$$P_{X,Y}(Y \in \bullet \mid X = x) = P_{X,Y}(Y \in \bullet \mid X = x, S = S(x)) = P_{X,Y}(Y \in \bullet \mid S = S(x)),$$

---

4. When  $Y$  is the prediction of  $X_{n+1}$  from  $\mathbf{X}_n$ , adequacy is also known as *predictive sufficiency* (Fortini et al., 2000) and, under certain conditions, it is equivalent to sufficiency and transitivity (Bahadur, 1954) or to total sufficiency (Lauritzen, 1974a). See Lauritzen (1974b) for a precise description of how these concepts are related.

which is equivalent to  $Y \perp\!\!\!\perp_{S(X)} X$ . Borrowing terminology from the graphical models literature, we say that  $S$  *d-separates*  $X$  and  $Y$  (Lauritzen, 1996). Therefore, adequacy is equivalent to sufficiency for  $\mathcal{P}_X$  and d-separation of  $X$  and  $Y$ , for all distributions in  $\mathcal{P}_{X,Y}$ .

### 3.2. Noise Outsourcing and Conditional Independence

Noise outsourcing is a standard technical tool from measure theoretic probability, where it is also known by other names such as *transfer* (Kallenberg, 2002). For any two random variables  $X$  and  $Y$  taking values in nice spaces (e.g., Borel spaces), noise outsourcing says that there exists a functional representation of samples from the conditional distribution  $P_{Y|X}$  in terms of  $X$  and independent noise:  $Y \stackrel{\text{a.s.}}{=} f(\eta, X)$ . As noted in Section 1, the relevant property of  $\eta$  is its independence from  $X$ , and the uniform distribution could be replaced by any other random variable taking values in a Borel space, for example a standard normal on  $\mathbb{R}$ , and the result would still hold, albeit with a different  $f$ .

Basic noise outsourcing can be refined in the presence of conditional independence. Let  $S : \mathcal{X} \rightarrow \mathcal{S}$  be a statistic such that  $Y$  and  $X$  are conditionally independent, given  $S(X)$ :  $Y \perp\!\!\!\perp_{S(X)} X$ . The following basic result, upon which the results in Section 4 rely, says that if there is a statistic  $S$  that d-separates  $X$  and  $Y$ , then it is possible to represent  $Y$  as a noise-outsourced function of  $S$ .

**Lemma 5** *Let  $X$  and  $Y$  be random variables with joint distribution  $P_{X,Y}$ . Let  $\mathcal{S}$  be a standard Borel space and  $S : \mathcal{X} \rightarrow \mathcal{S}$  a measurable map. Then  $S(X)$  d-separates  $X$  and  $Y$  if and only if there is a measurable function  $f : [0, 1] \times \mathcal{S} \rightarrow \mathcal{Y}$  such that*

$$(X, Y) \stackrel{\text{a.s.}}{=} (X, f(\eta, S(X))) \quad \text{where } \eta \sim \text{Unif}[0, 1] \quad \text{and } \eta \perp\!\!\!\perp X. \quad (13)$$

*In particular,  $Y = f(\eta, S(X))$  has distribution  $P_{Y|X}$ .*

The proof is a straightforward application of a more general result given in Appendix A. Note that in general,  $f$  is measurable but need not be differentiable or otherwise have desirable properties, although for modeling purposes it can be limited to functions belonging to a tractable class (e.g., differentiable, parameterized by a neural network). Note also that the identity map  $S(X) = X$  trivially d-separates  $X$  and  $Y$ , so that  $Y \stackrel{\text{a.s.}}{=} f(\eta, X)$ , which is standard noise outsourcing (e.g., Austin, 2015, Lem. 3.1).

To make the connections between ideas clear, we state the following corollary of Theorem 5, which is proved by checking the definition of adequacy.

**Corollary 6** *Let  $S : \mathcal{X} \rightarrow \mathcal{S}$  be a sufficient statistic for the model  $\mathcal{P}_X$ . Then  $S$  is an adequate statistic for  $\mathcal{P}_{X,Y}$  if and only if, for each  $P_{X,Y} \in \mathcal{P}_{X,Y}$ , there exists a corresponding measurable function  $f : [0, 1] \times \mathcal{S} \rightarrow \mathcal{Y}$  such that (13) holds.*

## 4. Functional Representations of Probabilistic Symmetries

Theorem 5 shows that if it is possible to find conditions on  $P_{X,Y}$  such that there is a  $\mathcal{G}$ -invariant d-separating statistic  $S$ , then there exists a  $\mathcal{G}$ -invariant functional representation of  $Y$  that depends on  $X$  only through  $S(X)$ . Furthermore, families of probability measures (i.e., statistical models) consisting of such distributions correspond to families of functions

with the desired invariance property. The main results of this section establish necessary and sufficient conditions on  $P_{X,Y}$  to guarantee  $\mathcal{G}$ -invariant (Theorem 7) and  $\mathcal{G}$ -equivariant (Theorem 9) functional representations of  $Y$ ; they also establish generic properties of the associated d-separating statistics. Together, they form a general program for identifying function classes and neural network architectures corresponding to symmetry conditions.

The results of this section imply that the two approaches in Section 2 are equivalent up to stochasticity (which is not trivial; see Section 5.3): symmetry can be incorporated into an explicit probabilistic model by considering invariant and equivariant distributions, or it can be incorporated into an implicit model by considering invariant and equivariant stochastic functions. The deterministic functional models in Section 2.1 are special cases of the latter: for every invariant noise-outsourced function  $f$  such that  $Y = f(\eta, X)$ , we have that

$$\mathbb{E}[Y \mid M(X)] = \int_{[0,1]} f(\eta, X) d\eta = f'(X),$$

such that  $f'(X)$  is also invariant (and similarly for equivariant functions).

#### 4.1. Invariant Conditional Distributions

To state the first main result, on the function representation of invariant conditional distributions, some definitions are required. For a group  $\mathcal{G}$  acting on a set  $\mathcal{X}$ , the *orbit* of any  $x \in \mathcal{X}$  is the set of elements in  $\mathcal{X}$  that can be generated by applying the elements of  $\mathcal{G}$ . It is denoted  $\mathcal{G} \cdot x = \{g \cdot x; g \in \mathcal{G}\}$ . The *stabilizer*, or isotropy subgroup, of  $x \in \mathcal{X}$  is the subgroup of  $\mathcal{G}$  that leaves  $x$  unchanged:  $\mathcal{G}_x = \{g \in \mathcal{G}; g \cdot x = x\}$ . An *invariant statistic*  $S : \mathcal{X} \rightarrow \mathcal{S}$  is a measurable map that satisfies  $S(x) = S(g \cdot x)$  for all  $g \in \mathcal{G}$  and  $x \in \mathcal{X}$ . A *maximal invariant statistic*, or maximal invariant, is an invariant statistic  $M : \mathcal{X} \rightarrow \mathcal{S}$  such that  $M(x_1) = M(x_2)$  implies  $x_2 = g \cdot x_1$  for some  $g \in \mathcal{G}$ ; equivalently,  $M$  takes a different constant value on each orbit. The orbits partition  $\mathcal{X}$  into equivalence classes, on each of which  $M$  takes a different value.

By definition, an invariant distribution  $P_X$  is constant on any particular orbit. Consider the conditional distribution  $P_X(X \mid M(X) = m) := P_{X|m}(X)$ . Conditioning on the maximal invariant taking a particular value is equivalent to conditioning on  $X$  being in a particular orbit; for invariant  $P_X$ ,  $P_{X|m}$  is zero outside the orbit on which  $M(X) = m$ , and “uniform” on the orbit, modulo fixed points (see Footnote 14 in Appendix B). Furthermore, for any  $Y$  such that  $(g \cdot X, Y) \stackrel{d}{=} (X, Y)$  for all  $g \in \mathcal{G}$ ,  $M(X)$  contains all relevant information for predicting  $Y$  from  $X$ ; that is,  $Y \perp\!\!\!\perp_{M(X)} X$ . Figure 2 illustrates the structure. These high-level ideas, which are made rigorous in Appendix B, lead to the following functional representation of invariant conditional distributions.

**Theorem 7** *Let  $X$  and  $Y$  be random elements of Borel spaces  $\mathcal{X}$  and  $\mathcal{Y}$ , respectively, and  $\mathcal{G}$  a compact group acting measurably on  $\mathcal{X}$ . Assume that  $P_X$  is  $\mathcal{G}$ -invariant, and pick a maximal invariant  $M : \mathcal{X} \rightarrow \mathcal{S}$ , with  $\mathcal{S}$  another Borel space. Then  $P_{Y|X}$  is  $\mathcal{G}$ -invariant if and only if there exists a measurable function  $f : [0, 1] \times \mathcal{S} \rightarrow \mathcal{Y}$  such that*

$$(X, Y) \stackrel{\text{a.s.}}{=} (X, f(\eta, M(X))) \quad \text{with } \eta \sim \text{Unif}[0, 1] \text{ and } \eta \perp\!\!\!\perp X. \quad (14)$$



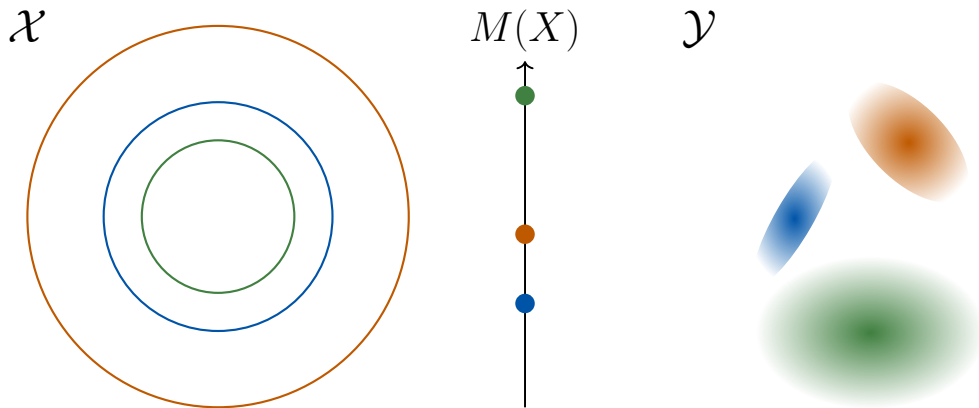


Figure 2: An illustration of structure of Theorem 7. The maximal invariant  $M$  (middle) acts as an index of the orbits of  $\mathcal{X}$  under  $\mathcal{G}$  (left), and mediates all dependence of  $Y$  on  $X$ , i.e.,  $P_{Y|X} = P_{Y|M(X)}$  (right). (Best viewed in color.)

The proof is given in Appendix B after the necessary intermediate technical results have been developed. Note that because  $f$  is a function of a maximal invariant, it is  $\mathcal{G}$ -invariant:  $f(\eta, M(g \cdot X)) = f(\eta, M(X))$ , almost surely. A maximal invariant always exists for sufficiently nice  $\mathcal{G}$  and  $\mathcal{X}$  (Hall et al., 1965); for example, define  $M : \mathcal{X} \rightarrow \mathbb{R}$  to be a function that takes a unique value on each orbit. Therefore,  $P_{X|M}$  can always be defined in such cases. The assumptions made in the present work, that  $\mathcal{G}$  acts measurably on the Borel space  $\mathcal{X}$ , allow for the existence of maximal invariants, and in many settings of interest a maximal invariant is straightforward to construct. For example, Sections 6 and 7 rely on constructing maximal invariants for applications of the results of this section to specific exchangeable structures.

Versions of Theorem 7 may hold for non-compact groups and more general topological spaces, but require considerably more technical details. At a high level,  $\mathcal{G}$ -invariant measures on  $\mathcal{X}$  may be disintegrated into product measures on suitable spaces  $\mathcal{S} \times \mathcal{Z}$  under fairly general conditions, though extra care is needed in order to ensure that such disintegrations exist. See, for example, Andersson (1982); Eaton (1989); Wijsman (1990); Schindler (2003); and especially Kallenberg (2017).

## 4.2. Equivariant Conditional Distributions

According to Theorem 7, *any* maximal invariant can be used to establish a functional representation of  $\mathcal{G}$ -invariant conditional distributions. If a particular type of equivariant function exists and is measurable, then it can be used to establish a functional representation of equivariant conditional distributions. Let  $\tau : \mathcal{X} \rightarrow \mathcal{G}$  be an equivariant function,

$$\tau(g \cdot x) = g \cdot \tau(x), \quad \text{for all } g \in \mathcal{G}, x \in \mathcal{X}. \quad (15)$$

For ease of notation, let  $\tau_x := \tau(x)$ , and denote by  $\tau_x^{-1}$  the inverse of  $\tau_x$  in  $\mathcal{G}$ , such that  $\tau_x^{-1} \cdot \tau_x = e$ .  $\tau$  has some remarkable properties that make it suitable for constructing equivariant functions.

**Lemma 8** *For a group  $\mathcal{G}$  acting measurably on Borel spaces  $\mathcal{X}$  and  $\mathcal{Y}$ , a representative equivariant  $\tau : \mathcal{X} \rightarrow \mathcal{G}$ , as defined in (15), has the following properties:*

- (i) *The function  $M_\tau : \mathcal{X} \rightarrow \mathcal{X}$  defined by  $M_\tau(x) = \tau_x^{-1} \cdot x$  is a maximal invariant.*
- (ii) *For any mapping  $b : \mathcal{X} \rightarrow \mathcal{Y}$ , the function*

$$f(x) = \tau_x \cdot b(\tau_x^{-1} \cdot x), \quad x \in \mathcal{X}, \quad (16)$$

*is  $\mathcal{G}$ -equivariant:  $f(g \cdot x) = g \cdot f(x)$ ,  $g \in \mathcal{G}$ .*

We call  $\tau$  a *representative equivariant* of its use in constructing the maximal invariant  $M_\tau(x) := \tau_x^{-1} \cdot x$ , which is a representative element from each orbit in  $\mathcal{X}$ . (We give examples in Sections 6 and 7.) The properties in Theorem 8 are used to establish the equivariant counterpart of Theorem 7. In essence, for  $\mathcal{G}$ -invariant  $P_X$ ,  $P_{Y|X}$  is conditionally  $\mathcal{G}$ -equivariant if and only if there exists a noise-outsourced function such that

$$g \cdot Y \stackrel{\text{a.s.}}{=} f(\eta, g \cdot X). \quad (17)$$

Observe that this is equivalent to  $f$  being  $\mathcal{G}$ -equivariant, as defined in (5), in the second argument:  $Y = e \cdot Y = f(\eta, e \cdot X) = f(\eta, X)$  and therefore

$$f(\eta, g \cdot X) = g \cdot Y = g \cdot f(\eta, X), \quad \text{a.s., for each } g \in \mathcal{G}.$$

It is straightforward to show that for  $X$  and  $Y$  constructed in this way, the resulting distributions have the desired equivariance property. The existence of an equivariant functional representation is harder to show, and the representative equivariant  $\tau$  plays a key role. We elaborate after stating the result.

**Theorem 9** *Let  $\mathcal{G}$  be a compact group acting measurably on Borel spaces  $\mathcal{X}$  and  $\mathcal{Y}$ , such that there exists a measurable representative equivariant  $\tau : \mathcal{X} \rightarrow \mathcal{G}$  satisfying (15). Suppose  $P_X$  is  $\mathcal{G}$ -invariant. Then  $P_{Y|X}$  is  $\mathcal{G}$ -equivariant if and only if there exists a measurable  $\mathcal{G}$ -equivariant function  $f : [0, 1] \times \mathcal{X} \rightarrow \mathcal{Y}$  satisfying (17) such that*

$$(X, Y) \stackrel{\text{a.s.}}{=} (X, f(\eta, X)) \quad \text{for each } g \in \mathcal{G}, \text{ with } \eta \sim \text{Unif}[0, 1] \text{ and } \eta \perp\!\!\!\perp X.$$

The proof of Theorem 9 adapts an argument due to Kallenberg (2005, Lem. 7.11). The proof, given in full in Appendix B.2, is rather technical but we sketch the basic ideas here.  $\mathcal{G}$ -invariance of  $P_X$  and  $\mathcal{G}$ -equivariance of  $P_{Y|X}$  is equivalent to the joint invariance  $(g \cdot X, g \cdot Y) \stackrel{\text{d}}{=} (X, Y)$ , for all  $g \in \mathcal{G}$ . As such, we require a frame of reference relative to which  $X$  and  $Y$  will equivary. The representative equivariant  $\tau$  plays this role. Recall that  $\tau$  is used to define a maximal invariant,  $M_\tau(x) := \tau_x^{-1} \cdot x$ , which maps  $x$  to the representative element of its orbit. The proof of Theorem 9 relies on establishing the

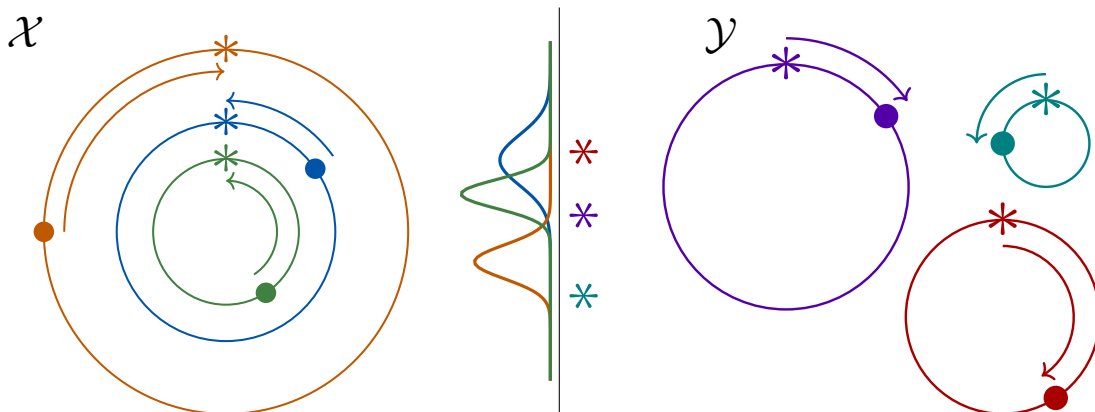


Figure 3: Illustration of the structure of Theorem 9. Each orbit of  $\mathcal{X}$  (left) induces a distribution over the orbits of  $\mathcal{Y}$  (middle). Given  $X$ , a sample  $Y$  is generated by: 1) obtaining the orbit representative  $\tau_X^{-1} \cdot X$ ; 2) sampling an orbit representative in  $\mathcal{Y}$ , conditioned on  $\tau_X^{-1} \cdot X$ ; and 3) applying  $\tau_X$  to the orbit representative in  $\mathcal{Y}$ . (Best viewed in color.)

conditional independence relationship  $\tau_X^{-1} \cdot Y \perp\!\!\!\perp_{M_\tau(X)} X$ . Noise outsourcing (Theorem 5) then implies that  $\tau_X^{-1} \cdot Y = b(\eta, M_\tau(X))$  for some  $b: [0, 1] \times \mathcal{X} \rightarrow \mathcal{Y}$ , and hence

$$Y = f(\eta, X) := \tau_X \cdot b(\eta, M_\tau(X)) = \tau_X \cdot b(\eta, \tau_X^{-1} \cdot X). \quad (18)$$

Finally, Theorem 8 (see Appendix B) establishes that  $f$  so defined is equivariant in the second argument:  $f(\eta, g \cdot X) = g \cdot f(\eta, X)$  for all  $g \in \mathcal{G}$ .

Observe that if the action of  $\mathcal{G}$  on  $\mathcal{Y}$  is trivial,  $\tau_X^{-1} \cdot Y \perp\!\!\!\perp_{M_\tau(X)} X$  reduces to  $Y \perp\!\!\!\perp_{M_\tau(X)} X$ , which is precisely the relationship needed to establish the invariant representation in Theorem 7. However, in the invariant case, *any* maximal invariant will d-separate  $X$  and  $Y$ , and therefore Theorem 7 is more general than simply applying Theorem 9 to the case with trivial group action on  $\mathcal{Y}$ . The additional assumptions in Theorem 9 account for the non-trivial action of  $\mathcal{G}$  on  $\mathcal{Y}$ , which requires setting up a fixed frame of reference through  $\tau$  and the orbit representatives  $M_\tau(X)$ .

Theorem 9 may hold for non-compact groups, but our proof for compact groups makes use of the normalized Haar measure; extending to non-compact groups requires additional technical overhead.

### 4.3. Symmetry-induced Adequate Statistics

The statistics upon which the functional representations in Theorems 7 and 9 rely capture all of the relevant structural information from  $X$ , and discard the rest. The utility of this orbit-resolving property can be characterized statistically. For suitably defined models—families of distributions satisfying the required invariance conditions—a maximal invariant is a sufficient statistic for the marginal model  $\mathcal{P}_X$ , and an adequate statistic for the joint

model  $\mathcal{P}_{X,Y}$ . Similarly, for models satisfying the equivariance conditions, a representative equivariant is an adequate statistic.

To make this precise, denote by  $\mathcal{P}_{X,Y}^{\text{inv}}$  the family of probability measures on  $\mathcal{X} \times \mathcal{Y}$  that satisfy  $(g \cdot X, Y) \stackrel{d}{=} (X, Y)$ , for all  $g \in \mathcal{G}$ . Similarly, denote by  $\mathcal{P}_{X,Y}^{\text{eq}}$  the family of probability measures on  $\mathcal{X} \times \mathcal{Y}$  that satisfy  $(g \cdot X, g \cdot Y) \stackrel{d}{=} (X, Y)$ , for all  $g \in \mathcal{G}$ . Let  $\mathcal{P}_{X, \tau_X^{-1} \cdot Y}$  denote the family of probability measures obtained from  $\mathcal{P}_{X,Y}^{\text{eq}}$  via  $P_{X,Y} \mapsto P_X P_{Y|X} \circ \tau_X^{-1}$ .

**Theorem 10** *Let  $\mathcal{G}$  be a compact group acting measurably on standard Borel spaces  $\mathcal{X}$  and  $\mathcal{Y}$ , and let  $\mathcal{S}$  be another Borel space. Then:*

- (i) *Any maximal invariant  $M : \mathcal{X} \rightarrow \mathcal{S}$  on  $\mathcal{X}$  under  $\mathcal{G}$  is an adequate statistic of  $X$  for  $Y$  with respect to any model  $\mathcal{P}_{X,Y} \subseteq \mathcal{P}_{X,Y}^{\text{inv}}$ . In particular, to each  $P_{X,Y} \in \mathcal{P}_{X,Y}^{\text{inv}}$  corresponds a measurable  $\mathcal{G}$ -invariant function  $f : [0, 1] \times \mathcal{S} \rightarrow \mathcal{Y}$  as in (14).*
- (ii)  *$M_\tau$  as in Theorem 8 is an adequate statistic of  $X$  for  $\tau_X^{-1} \cdot Y$  with respect to any model  $\mathcal{P}_{X,Y} \subseteq \mathcal{P}_{X, \tau_X^{-1} \cdot Y}$ , and to each  $P_{X,Y} \in \mathcal{P}_{X, \tau_X^{-1} \cdot Y}$  there corresponds a measurable function  $b : [0, 1] \times \mathcal{X} \rightarrow \mathcal{Y}$  as in (18) such that  $\tau_X \cdot b(\eta, \tau_X^{-1} \cdot X)$  is  $\mathcal{G}$ -equivariant.*

Theorem 10 puts explicit probability models in correspondence with implicit functional models. It is a consequence of Theorem 20 in Appendix B, which establishes the necessary structural properties of all invariant distributions. To illustrate, consider again the example of an exchangeable  $\mathcal{X}$ -valued sequence.

**Example 3 (Adequate statistics of an exchangeable sequence)** *Let  $\mathbf{X}_n$  be an exchangeable  $\mathcal{X}$ -valued sequence.  $\mathbb{S}_n$  acts on any point  $\mathbf{x}_n \in \mathcal{X}^n$  by permuting its indices, which defines an action on  $\mathcal{X}^n$ ; the orbit of  $\mathbf{x}_n$  is the set of sequences that can be obtained from  $\mathbf{x}_n$  by applying a permutation. The empirical measure is a well-known sufficient statistic for  $\mathcal{P}_{\mathbf{X}_n}^{\mathbb{S}_n}$ ; it is easy to see that the empirical measure is also a maximal invariant (see Section 6.1). If  $\mathcal{X} = \mathbb{R}$  (or any other set with a total order), then so too is the vector of order statistics,  $\mathbf{x}_n^\uparrow = (x_{(1)}, \dots, x_{(n)})$ , with  $x_{(1)} \leq \dots \leq x_{(n)}$ . A representative equivariant may be obtained by defining  $\tau_{\mathbf{x}_n}^{-1}$  as any permutation  $\pi$  (not necessarily unique) that satisfies  $\pi \cdot \mathbf{x}_n = \mathbf{x}_n^\uparrow$ .*

*Exchangeable random structures.* Sections 6 and 7 work out the details of the program outlined in Section 5.1 to various exchangeable random structures, and relates the resulting representations to some recent (Zaheer et al., 2017; Gilmer et al., 2017; Hartford et al., 2018) and less recent (Shawe-Taylor, 1989) work in the neural networks literature. Exchangeability is distributional invariance under the action of  $\mathbb{S}_n$  (or other groups defined by composing  $\mathbb{S}_n$  in different ways for other data structures). Example 3 states that the empirical measure is a maximal invariant of  $\mathbb{S}_n$  acting on  $\mathcal{X}^n$ ; suitably defined generalizations of the empirical measure are maximal invariants for other exchangeable structures. With these maximal invariants, obtaining a functional representation of an invariant conditional distribution is straightforward.

With an additional conditional independence assumption that is satisfied by most neural network architectures, Theorem 9 can be refined to obtain a detailed functional representation of the relationship between  $\mathbb{S}_n$ -equivariant random variables (Theorem 12). That

refinement relies on the particular subgroup structure of  $\mathbb{S}_n$ , and it raises the question, not pursued here, of whether, and under what conditions, other groups may yield similar refinements.

## 5. Practical Implications and Considerations

In order to make the previous sections' theory more useful to practitioners, we formulate a general program for model specification, briefly discuss computational considerations, and interpret some aspects of the theory in the context of the recent literature.

### 5.1. A Program for Obtaining Symmetric Functional Representations

The previous section suggests that for a model consisting of  $\mathcal{G}$ -invariant distributions, if an adequate statistic can be found then all distributions in the conditional model,  $P_{Y|X} \in \mathcal{P}_{Y|X}$  have a noise-outsourced functional representation in terms of the adequate statistic, as in (13). Furthermore, Theorem 10 says that any maximal invariant is an adequate statistic under such a model. Therefore, for a specified compact group  $\mathcal{G}$  acting measurably on  $\mathcal{X}$  and  $\mathcal{Y}$ , a program for functional model specification through distributional symmetry is as follows:

PROGRAM FOR  $\mathcal{G}$ -EQUIVARIANT AND -INVARIANT MODEL SPECIFICATION.

- (i) Determine a maximal  $\mathcal{G}$ -invariant statistic  $M : \mathcal{X} \rightarrow \mathcal{S}$ .
- (ii) Fix a representative equivariant under  $\mathcal{G}$ ,  $\tau : \mathcal{X} \rightarrow \mathcal{G}$ , that can be used to map any element of  $\mathcal{X}$  to its orbit's representative.
- (iii) Specify the model with a function class  $\mathcal{F}_{\tau, M}$ , consisting of compositions of noise-outsourced equivariant functions and a final noise-outsourced invariant function.

### 5.2. Computing Maximal Invariants and Representative Equivariants

In theory, a maximal invariant can be computed by specifying a representative element for each orbit. In practice, this can always be done when  $\mathcal{G}$  is discrete because the relevant elements can be enumerated and reduced to permutation operations. Several systems for computational discrete mathematics, such as Mathematica (Wolfram Research, Inc., 2018) and GAP (The GAP Group, 2018), have built-in functions to do so. For continuous groups it may not be clear how to compute a maximal invariant, and it may be impractical. Furthermore, maximal invariants are not unique, and some may be better suited to a particular application than others. As such, they are best handled on a case-by-case basis, depending on the problem at hand. Some examples from the classical statistics literature are reviewed in Lehmann and Romano (2005, Ch. 6) and Eaton (1989, Ch. 2); Kallenberg (2017, Ch. 7) presents a generic method based on so-called projection maps. Sections 6 and 7 apply the theory of this section to exchangeable structures by explicitly constructing maximal invariants.

Likewise, a representative equivariant  $\tau$ , as used in Theorem 9, always exists for a discrete group acting on a Borel space  $\mathcal{X}$ , a proof of which is given in Kallenberg (2005, Lem. 7.10). If  $\mathcal{G}$  is compact and acts measurably on Borel spaces  $\mathcal{X}$  and  $\mathcal{Y}$ , then a representative

equivariant exists; see Schindler (2003, Remark 2.46). In more general settings,  $\tau$  with the desired properties may not be well-defined or it may not be measurable, in which case a  $\mathcal{G}$ -invariant probability kernel (called an *inversion kernel*) may be constructed from a maximal invariant to obtain similar representations; see Kallenberg (2017, Ch. 7). Whether a representative equivariant is readily computed is a different matter. If, as in Example 3, a maximal invariant corresponding to a representative element of each orbit is computed, it corresponds to  $M_\tau$  and the representative equivariant is computed as a by-product. However, even this can be difficult. For example, in the case of graphs (Section 7), finding a maximal invariant statistic is equivalent to finding a complete graph invariant (e.g., Lovász, 2012); finding a representative equivariant is equivalent to graph canonization, which is at least as computationally hard as the graph isomorphism problem.

### 5.3. On the Advantages of Stochasticity

Much of deep learning (and machine learning in general) is conducted in the framework of learning a deterministic function that does not have stochastic component apart from the data. Although that approach has been successful in a wide variety of settings, and the methods are often more accessible to users without training in probability, there are some fundamental limitations to the approach. One body of work has argued that injecting noise into a neural network during training has beneficial effects on training (e.g., Srivastava et al., 2014; Wager et al., 2014).

We focus here on a different aspect. Let  $X$  and  $Y \in \mathbb{R}$  be the input and output data, respectively, and let  $Y' = f(X)$ , for some  $f : \mathcal{X} \rightarrow \mathbb{R}$ . A basic result about conditional expectation (see, e.g., Çinlar, 2011, Thm. IV.1.17) says that an  $\mathbb{R}$ -valued random variable  $Y'$  is a deterministic (measurable) function of  $X$  if and only if  $Y'$  corresponds to the conditional expectation of some (typically not unique) random variable (say  $Z$ ) given  $X$ . That is,

$$Y' = f(X) \iff Y' = \mathbb{E}[Z | X].$$

Much of machine learning relies on (approximately) minimizing a risk  $\mathbb{E}[\ell(f(X), Y)]$  with respect to  $f$ . The best such a procedure can do is to learn the conditional expectation  $\mathbb{E}[Y | X]$ ; the model does not allow  $Y$  to exhibit any randomness apart from that in  $X$ . Except for certain cases (see below), important information may be lost.

*Energy-based models.* The extension from conditional expectation  $\mathbb{E}[Y | X]$  to conditional distribution  $P_{Y|X}$  is typically not unique. One principled way of bridging the gap is via maximum entropy, which stipulates that  $P_{Y|X}$  should be the distribution with the maximum entropy,  $H(P_{Y|X})$ , subject to a constraint on the conditional expectation, e.g.,  $\mathbb{E}[Y | X] = s$ . The result is a distribution with density

$$p_{Y|X}(y) \propto \exp(-\beta y) = \arg \min_{p_{Y|X}} \mathcal{L}(p_{Y|X}, \beta) = -H(p_{Y|X}) + \beta(\mathbb{E}_{p_{Y|X}}[Y | X] - s),$$

that minimizes the free energy  $\mathcal{L}(p_{Y|X}, \beta)$ . This approach, which has a long history in statistical physics (macrocanonical models) and classical statistics (exponential families) (Jaynes, 1957), was the starting point for work on invariant wavelet-based models (Bruna and Mallat, 2013, 2018; Gao et al., 2019), which can be described in terms of invariant sufficient statistics. More broadly, this approach fits in the framework of energy-based learning

with Gibbs distributions (LeCun et al., 2006; Grathwohl et al., 2020), which subsumes much of the standard machine learning practice of training deterministic functions with (regularized) empirical risk minimization. Despite the long history and easy of interpretation of energy-based models, more sophisticated approaches that incorporate randomness in the network itself seem attractive, for the reasons detailed above.

*Simpler proofs and clearer structure.* From a theoretical perspective, the probabilistic approach often yields simpler proofs of general results than their deterministic counterparts. Although this assessment can be somewhat subjective, this paper is an example: most of our results rely on establishing conditional independence and appealing to Theorem 5. This includes Theorem 11, which characterizes the representation of permutation-invariant distributions. Its the proof of its deterministic counterpart by Zaheer et al. (2017), relies on rather lengthy arguments about symmetric polynomials. Furthermore, the structure driving the representation becomes, in our opinion, much clearer from the perspective of sufficiency.

#### 5.4. Choosing a Function Class

Theorems 7 and 9 are very general existence results about the representations of symmetric distributions, and in that sense they are related to the recent literature on *universal approximations* (as in Cybenko, 1989; Hornik, 1991) of invariant and equivariant functions (e.g., Yarotsky, 2018; Maron et al., 2019; Keriven and Peyré, 2019; Segol and Lipman, 2020; Ravanbakhsh, 2020). Those results rely on approximating to arbitrary precision  $\mathcal{G}$ -symmetric polynomials. The most practically useful work considers the universal approximation problem in terms of composing linear maps with pointwise activations (Maron et al., 2019; Segol and Lipman, 2020; Ravanbakhsh, 2020). For finite groups, that approach also gives upper and lower bounds on the degree of the polynomial required. Despite being conceptually related to that work, Theorems 7 and 9 are *exact* representation results for  $\mathcal{G}$ -invariant and -equivariant distributions. Universal distributional approximation is an open question.

However, like the universal approximation literature, Theorems 7 and 9 are mostly silent about the practical matter of *which* class of functions should be used for a model. Detailed analysis of this type of problem, which is beyond the scope of this paper, is an active area of research requiring detailed analysis of the trade-off between model complexity, computation, and problem-specific requirements. Bietti and Mairal (2019) takes some initial steps in this direction.

*Convolutions and beyond.* The most widely used equivariant layers correspond to convolutions (with respect to  $\mathcal{G}$ ) with an equivariant convolution kernel. Recent work shows that they are the *only* equivariant linear maps (Kondor and Trivedi, 2018; Cohen et al., 2019). The details of this program have been implemented for a growing number of groups; see Cohen et al. (2019, Appendix D) for a comprehensive list. Modifications of the basic convolutional filters have appeared in the form of convolutional capsules (e.g., Sabour et al., 2017; Hinton et al., 2018; Lenssen et al., 2018), and convolution was recently combined with attention mechanisms (Bello et al., 2019; Romero et al., 2020). The representation result in Theorem 9 says nothing about linearity; non-convolutional equivariant architectures based on self-attention (Parmar et al., 2019) have recently appeared.

## 6. Learning from Finitely Exchangeable Sequences

In this section,<sup>5</sup> the program described in Section 5.1 is fully developed for the case where the conditional distribution of  $Y$  given  $X$  has invariance or equivariance properties with respect to  $\mathbb{S}_n$ . Deterministic examples have appeared in the neural networks literature (Shawe-Taylor, 1989; Zaheer et al., 2017; Ravanbakhsh et al., 2017; Murphy et al., 2019), and the theory developed in here establishes the necessary and sufficient functional forms for permutation invariance and equivariance, in both the stochastic and deterministic cases.

Throughout this section, the input  $X = \mathbf{X}_n$  is a sequence of length  $n$ , and  $Y$  is an output variable, whose conditional distribution given  $\mathbf{X}_n$  is to be modeled.<sup>6</sup> Recall from Section 2.2 that  $P_{Y|\mathbf{X}_n}$  is  $\mathbb{S}_n$ -invariant if  $Y|\mathbf{X}_n \stackrel{d}{=} Y|\pi \cdot \mathbf{X}_n$  for each permutation  $\pi \in \mathbb{S}_n$  of the input sequence. Alternatively, if  $Y = \mathbf{Y}_n$  is also a sequence of length  $n$ ,<sup>7</sup> then we say that  $\mathbf{Y}_n$  given  $\mathbf{X}_n$  is  $\mathbb{S}_n$ -equivariant if  $\mathbf{Y}_n|\mathbf{X}_n \stackrel{d}{=} \pi \cdot \mathbf{Y}_n|\pi \cdot \mathbf{X}_n$  for all  $\pi \in \mathbb{S}_n$ . In both cases these symmetry properties stem from the assumption that the ordering in the input sequence  $\mathbf{X}_n$  does not matter; that is, the distribution of  $\mathbf{X}_n$  is finitely exchangeable:  $\mathbf{X}_n \stackrel{d}{=} \pi \cdot \mathbf{X}_n$  for each  $\pi \in \mathbb{S}_n$ . Recall that  $P_{\mathbf{X}_n}$  and  $P_{Y|\mathbf{X}_n}$  denote the marginal and conditional distributions respectively,  $\mathcal{P}_{\mathbf{X}_n}^{\mathbb{S}_n}$  is the family of distributions on  $\mathcal{X}^n$  that are  $\mathbb{S}_n$ -invariant, and  $\mathcal{P}_{Y|\mathbf{X}_n}^{\mathbb{S}_n}$  is the family of conditional distributions on  $\mathcal{Y}$  given  $\mathbf{X}_n$  that are  $\mathbb{S}_n$ -invariant.

The primary conceptual matter is the central role of the empirical measure  $\mathbb{M}_{\mathbf{X}_n}$ . Exchangeable sequences have been studied in great detail, and the sufficiency of the empirical measure for  $\mathcal{P}_{\mathbf{X}_n}^{\mathbb{S}_n}$  is well-known (e.g., Diaconis and Freedman (1980a); Kallenberg (2005, Prop. 1.8)). It is also straightforward to show the adequacy of the empirical measure for  $\mathcal{P}_{\mathbf{X}_n, Y}^{\mathbb{S}_n}$  using methods that are not explicitly group theoretic. Alternatively, it is enough to show that the empirical measure is a maximal invariant of  $\mathcal{X}^n$  under  $\mathbb{S}_n$  and then apply Theorem 10. In either case, the results of the previous sections imply a noise-outsourced functional representation of  $\mathbb{S}_n$ -invariant conditional distributions (Section 6.1). The previous sections also imply a representation for  $\mathbb{S}_n$ -equivariant conditional distributions, but under an additional conditional independence assumption a more detailed representation can be obtained due to the structure of  $\mathbb{S}_n$  (Section 6.2).

### 6.1. $\mathbb{S}_n$ -Invariant Conditional Distributions

The following special case of Theorems 7 and 10 establishes a functional representation for all  $\mathbb{S}_n$ -invariant conditional distributions.

**Theorem 11** *Let  $\mathbf{X}_n \in \mathcal{X}^n$ , for some  $n \in \mathbb{N}$ , be an exchangeable sequence, and  $Y \in \mathcal{Y}$  some other random variable. Then  $P_{Y|\mathbf{X}_n}$  is  $\mathbb{S}_n$ -invariant if and only if there is a measurable function  $f : [0, 1] \times \mathcal{M}(\mathcal{X}) \rightarrow \mathcal{Y}$  such that*

$$(\mathbf{X}_n, Y) \stackrel{\text{a.s.}}{=} (\mathbf{X}_n, f(\eta, \mathbb{M}_{\mathbf{X}_n})) \quad \text{where } \eta \sim \text{Unif}[0, 1] \quad \text{and } \eta \perp\!\!\!\perp \mathbf{X}_n. \quad (19)$$

*Furthermore,  $\mathbb{M}_{\mathbf{X}_n}$  is sufficient for the family  $\mathcal{P}_{\mathbf{X}_n}^{\mathbb{S}_n}$ , and adequate for the family  $\mathcal{P}_{\mathbf{X}_n, Y}^{\mathbb{S}_n}$ .*

---

5. The main results of this section appeared in an extended abstract (Bloem-Reddy and Teh, 2018).  
6. Note that  $\mathbf{X}_n$  may represent a set (i.e., there are no repeated values) or a multi-set (there may be repeated values). It depends entirely on  $P_{\mathbf{X}_n}$ : if  $P_{\mathbf{X}_n}$  places all of its probability mass on sequences  $x_n \in \mathcal{X}^n$  that do not have repeated values, then  $\mathbf{X}_n$  represents a set almost surely. Otherwise,  $\mathbf{X}_n$  represents a multi-set. The results of this section hold in either case.  
7. In general,  $\mathbf{Y}$  need not be of length  $n$ , but the results are much simpler when it is; see Section 6.4.



**Proof** With Theorem 7, we need only prove that  $\mathbb{M}_{\mathbf{X}_n}$  is a maximal invariant of  $\mathcal{X}^n$  under  $\mathbb{S}_n$ . Clearly,  $\mathbb{M}_{\mathbf{X}_n}$  is  $\mathbb{S}_n$ -invariant. Now, let  $\mathbf{X}'_n$  be another sequence such that  $\mathbb{M}_{\mathbf{X}'_n} = \mathbb{M}_{\mathbf{X}_n}$ . Then  $\mathbf{X}'_n$  and  $\mathbf{X}_n$  contain the same elements of  $\mathcal{X}$ , and therefore  $\mathbf{X}'_n = \pi \cdot \mathbf{X}_n$  for some  $\pi \in \mathbb{S}_n$ , so  $\mathbb{M}_{\mathbf{X}_n}$  is a maximal invariant. ■

The second claim follow from Theorem 10. ■

*Modeling  $\mathbb{S}_n$ -invariance with neural networks.* Theorem 11 is a general characterization of  $\mathbb{S}_n$ -invariant conditional distributions. It says that all such conditional distributions must have a noise-outsourced functional representation given by  $Y = f(\eta, \mathbb{M}_{\mathbf{X}_n})$ . Recall that  $\mathbb{M}_{\mathbf{X}_n} = \sum_{i=1}^n \delta_{X_i}$ . An atom  $\delta_X$  can be thought of as a measure-valued generalization of a one-hot encoding to arbitrary measurable spaces, so that  $\mathbb{M}_{\mathbf{X}_n}$  is a sum-pooling of encodings of the inputs (which removes information about the ordering of  $\mathbf{X}_n$ ), and the output  $Y$  is obtained by passing that, along with independent outsourced noise  $\eta$ , through a function  $f$ . In case the conditional distribution is deterministic, the outsourced noise is unnecessary, and we simply have  $Y = f(\mathbb{M}_{\mathbf{X}_n})$ .

From a modeling perspective, one choice for (stochastic) neural network architectures that are  $\mathbb{S}_n$ -invariant is

$$Y = f\left(\eta, \sum_{i=1}^n \phi(X_i)\right), \quad (20)$$

where  $f$  and  $\phi$  are arbitrary neural network modules, with  $\phi$  interpreted as an embedding function of input elements into a high-dimensional space (see first panel of Fig. 4). These embeddings are sum-pooled, and passed through a second neural network module  $f$ . This architecture can be made to approximate arbitrarily well any  $\mathbb{S}_n$ -invariant conditional distribution (c.f., Hornik et al., 1989). Roughly,  $\phi(X)$  can be made arbitrarily close to a one-hot encoding of  $X$ , which can in turn be made arbitrarily close to an atom  $\delta_X$  by increasing its dimensionality, and similarly the neural module  $f$  can be made arbitrarily close to any desired function. Below, we revisit an earlier example and give some new ones.

**Example 4 (Deep Sets:  $\mathbb{S}_n$ -invariant functions of sequences, revisited)** *The architecture derived above is exactly the one described in Example 1. Theorem 11 generalizes the result in Zaheer et al. (2017), from deterministic functions to conditional distributions. The proof technique is also significantly simpler and sheds light on the core concepts underlying the functional representations of permutation invariance.*<sup>8</sup>

In general, a function of  $\mathbb{M}_{\mathbf{X}_n}$  is a function of  $\mathbf{X}_n$  that discards the order of its elements. That is, functions of  $\mathbb{M}_{\mathbf{X}_n}$  are permutation-invariant functions of  $\mathbf{X}_n$ . The sum-pooling in (20) gives rise to one such class of functions. Other permutation invariant pooling operations can be used. For example: product, maximum, minimum, log-sum-exp, mean, median, and

---

8. The result in Zaheer et al. (2017) holds for sets of arbitrary size when  $\mathcal{X}$  is countable and for fixed size when  $\mathcal{X}$  is uncountable. We note that the same is true for Theorem 11, for measure-theoretic reasons: in the countable  $\mathcal{X}$  case, the power sets of  $\mathbb{N}$  form a valid Borel  $\sigma$ -algebra; for uncountable  $\mathcal{X}$ , e.g.,  $\mathcal{X} = \mathbb{R}$ , there may be non-Borel sets and therefore the power sets do not form a Borel  $\sigma$ -algebra on which to define a probability distribution using standard techniques.

percentiles have been used in various neural network architectures. Any such function can be written in the form  $f(\eta, \mathbb{M}_{\mathbf{X}_n})$ , by absorbing the pooling operation into  $f$  itself. The following examples illustrate.

**Example 5 (Pooling using Abelian groups or semigroups)** *A group  $\mathcal{G}$ , with binary operator  $\oplus$ , is Abelian if its elements commute:  $g \oplus h = h \oplus g$  for all  $g, h \in \mathcal{G}$ . Examples are  $(\mathbb{R}_+, \times)$  and  $(\mathbb{Z}, +)$ . A semigroup is a group without the requirements for inverse and identity elements. Examples are  $(\mathbb{R}_+, \vee)$  ( $\vee$  denotes maximum:  $x \vee y = \max\{x, y\}$ ) and  $(\mathbb{R}_+, \wedge)$  ( $\wedge$  denotes minimum). For an Abelian group or semigroup  $\mathcal{G}$  and a map  $\phi : \mathcal{X} \rightarrow \mathcal{G}$ , a  $\mathbb{S}_n$ -invariant conditional distribution of  $Y$  given a sequence  $\mathbf{X}_n$  can be constructed with  $f : [0, 1] \times \mathcal{G} \rightarrow \mathcal{Y}$  as*

$$Y = f(\eta, \phi(X_1) \oplus \cdots \oplus \phi(X_n)) . \tag{21}$$

**Example 6 (Pooling using U-statistics)** *Given a permutation invariant function of  $k \leq n$  elements,  $\phi_k : \mathcal{X}^k \rightarrow \mathcal{S}$ , a permutation invariant conditional distribution can be constructed with  $f : [0, 1] \times \mathcal{S} \rightarrow \mathcal{Y}$  as*

$$f\left(\eta, \binom{n}{k}^{-1} \sum_{\{i_1, \dots, i_k\} \in [n]} \phi_k(X_{i_1}, \dots, X_{i_k})\right) \tag{22}$$

where the pooling involves averaging over all  $k$ -element subsets of  $[n]$ . The average is a  $U$ -statistic (e.g., Cox and Hinkley, 1974), and examples include the sample mean ( $k = 1$  and  $\phi_k(x) = x$ ), the sample variance ( $k = 2$  and  $\phi_k(x, y) = \frac{1}{2}(x - y)^2$ ), and estimators of higher-order moments, mixed moments, and cumulants.

Murphy et al. (2019) developed a host of generalizations to the basic first-order pooling functions from Example 1 (Deep Sets:  $\mathbb{S}_n$ -invariant functions of sequences), many of them corresponding to  $k$ -order  $U$ -statistics, and developed tractable computational techniques that approximate the average over  $k$ -element subsets by random sampling of permutations.

Exchangeability plays a central role in a growing body of work in the deep learning literature, particularly when deep learning methods are combined with Bayesian ideas. Examples include Edwards and Storkey (2017); Garnelo et al. (2018); Korshunova et al. (2018), and the following.

**Example 7 (Neural networks for exchangeable genetic data)** *In work by Chan et al. (2018),  $X_i \in \{0, 1\}^d$  is a binary  $d$ -dimensional vector indicating the presence or absence of  $d$  single nucleotide polymorphisms in individual  $i$ . The individuals are treated as being exchangeable, forming an exchangeable sequence  $\mathbf{X}_n$  of  $\{0, 1\}^d$ -valued random variables. Chan et al. (2018) analyze the data using a neural network where each vector  $X_i$  is embedded into  $\mathbb{R}^d$  using a convolutional network, the pooling operation is the element-wise mean of the top decile, and the final function is parameterized by a fully-connected network. They demonstrated empirically that encoding permutation invariance into the network architecture led to faster training and higher test set accuracy.*

The final example, from the statistical relational artificial intelligence literature, uses the sufficiency of the empirical measure to characterize the complexity of inference algorithms for exchangeable sequences.

**Example 8 (Lifted inference)** *Niepert and Van den Broeck (2014) studied the tractability of exact inference procedures for exchangeable models, through so-called lifted inference. One of their main results shows that if  $\mathbf{X}_n$  is a finitely exchangeable sequence of  $\mathcal{X}^d$ -valued random variables on a discrete domain (i.e., each element of  $\mathbf{X}_n$  is a  $d$ -dimensional vector of discrete random variables) then there is a sufficient statistic  $S$ , and probabilistic inference (defined as computing marginal and conditional distributions) based on  $S$  has computational complexity that is polynomial in  $d \times n$ . In the simplest case, where  $\mathcal{X} = \{0, 1\}$ ,  $S$  is constructed as follows: encode all possible  $d$ -length binary vectors with unique bit strings  $b_k \in \{0, 1\}^d$ ,  $k \in [2^d]$ , and let  $S(\mathbf{X}_n) = (c_1, \dots, c_{2^d})$  where  $c_k = \sum_{i=1}^n \delta_{X_i}(b_k)$ . Although not called the empirical measure by the authors,  $S$  is precisely that.*

## 6.2. $\mathbb{S}_n$ -Equivariant Conditional Distributions

Let  $\mathbf{X}_n$  be an input sequence of length  $n$ , and  $\mathbf{Y}_n$  an output sequence. Theorem 9 shows that if the conditional distribution of  $\mathbf{Y}_n$  given  $\mathbf{X}_n$  is  $\mathbb{S}_n$ -equivariant, then  $\mathbf{Y}_n$  can be expressed in terms of a noisy  $\mathbb{S}_n$ -equivariant function of  $\mathbf{X}_n$ . If the elements of  $\mathbf{Y}_n$  are assumed to be conditionally independent given  $\mathbf{X}_n$ , then by using properties of the finite symmetric group, we obtain a more detailed representation of  $\mathbf{Y}_n$  conditioned on  $\mathbf{X}_n$ . (The implications of the conditional independence assumption are discussed below.)

The resulting theorem is a one-dimensional special case of a more general representation theorem for exchangeable  $d$ -dimensional arrays (Theorem 23 in Appendix C). The following simplified proof for sequences shows how the necessary conditional independence relationships are established and provides a template for proving the more general result. The  $d = 2$  case, which corresponds to graphs and networks, is taken up in Section 7.

**Theorem 12** *Let  $\mathbf{X}_n \in \mathcal{X}^n$  be an exchangeable sequence and  $\mathbf{Y}_n \in \mathcal{Y}^n$  another random sequence, and assume that  $Y_i \perp\!\!\!\perp_{\mathbf{X}_n} (\mathbf{Y}_n \setminus Y_i)$ , for each  $i \in [n]$ . Then  $P_{\mathbf{Y}_n | \mathbf{X}_n}$  is  $\mathbb{S}_n$ -equivariant if and only if there is a measurable function  $f : [0, 1] \times \mathcal{X} \times \mathcal{M}(\mathcal{X}) \rightarrow \mathcal{Y}$  such that*

$$(\mathbf{X}_n, \mathbf{Y}_n) \stackrel{\text{a.s.}}{=} (\mathbf{X}_n, (f(\eta_i, X_i, \mathbb{M}_{\mathbf{X}_n}))_{i \in [n]}) \quad \text{where} \quad \eta_i \stackrel{\text{i.i.d.}}{\sim} \text{Unif}[0, 1] \quad \text{and} \quad \eta_i \perp\!\!\!\perp \mathbf{X}_n. \quad (23)$$

**Proof** For the forward direction, suppose  $\mathbf{Y}_n$  is conditionally  $\mathbb{S}_n$ -equivariant given  $\mathbf{X}_n$ . For a fixed  $i \in [n]$ , let  $\mathbf{X}_{n \setminus i} := \mathbf{X}_n \setminus X_i$  be  $\mathbf{X}_n$  with its  $i$ th element removed, and likewise for  $\mathbf{Y}_{n \setminus i}$ . The proof in this direction requires that we establish the conditional independence relationship

$$Y_i \perp\!\!\!\perp_{(X_i, \mathbb{M}_{\mathbf{X}_n})} (\mathbf{X}_n, \mathbf{Y}_{n \setminus i}), \quad (24)$$

and then apply Theorem 5.

To that end, let  $\mathbb{S}_{n \setminus i}$  be the stabilizer of  $i$ , i.e., the subgroup of  $\mathbb{S}_n$  that fixes element  $i$ .  $\mathbb{S}_{n \setminus i}$  consists of permutations  $\pi_{\setminus i} \in \mathbb{S}_n$  for which  $\pi_{\setminus i}(i) = i$ . The action of  $\pi_{\setminus i}$  on  $\mathbf{X}_n$  fixes  $X_i$ ; likewise it fixes  $Y_i$  in  $\mathbf{Y}_n$ . By Theorem 1,  $(\mathbf{X}_n, \mathbf{Y}_n) \stackrel{\text{d.}}{=} (\pi_{\setminus i} \cdot \mathbf{X}_n, \pi_{\setminus i} \cdot \mathbf{Y}_n)$ , so that, marginalizing out  $\mathbf{Y}_{n \setminus i}$  yields  $(\mathbf{X}_n, Y_i) \stackrel{\text{d.}}{=} (\pi_{\setminus i} \cdot \mathbf{X}_n, Y_i)$  for each  $\pi_{\setminus i} \in \mathbb{S}_{n \setminus i}$ . Moreover,  $\mathbb{S}_{n \setminus i}$  forms a subgroup and is homomorphic to  $\mathbb{S}_{n-1}$ , so that the previous distributional equality is equivalent to

$$(\mathbf{X}_{n \setminus i}, (X_i, Y_i)) \stackrel{\text{d.}}{=} (\pi' \cdot \mathbf{X}_{n \setminus i}, (X_i, Y_i)) \quad \text{for each} \quad \pi' \in \mathbb{S}_{n-1}.$$

Theorem 11, with input  $\mathbf{X}_{n \setminus i}$  and output  $(X_i, Y_i)$  then implies  $(X_i, Y_i) \perp\!\!\!\perp_{\mathbb{M}_{\mathbf{X}_{n \setminus i}}} \mathbf{X}_{n \setminus i}$ . Conditioning on  $X_i$  as well gives  $Y_i \perp\!\!\!\perp_{(X_i, \mathbb{M}_{\mathbf{X}_{n \setminus i}})} \mathbf{X}_n$ , marginally for each  $Y_i$ . With the assumption of mutual conditional independence among  $\mathbf{Y}_n$  conditioned on  $\mathbf{X}_n$ , the marginal conditional independence also holds jointly, and by the chain rule for conditional independence (Kallenberg, 2002, Prop. 6.8),

$$Y_i \perp\!\!\!\perp_{(X_i, \mathbb{M}_{\mathbf{X}_{n \setminus i}})} (\mathbf{X}_n, \mathbf{Y}_{n \setminus i}) . \quad (25)$$

Because conditioning on  $(X_i, \mathbb{M}_{\mathbf{X}_{n \setminus i}})$  is the same as conditioning on  $(X_i, \mathbb{M}_{\mathbf{X}_n})$ , (25) is equivalent to the key conditional independence relationship (24).

By Theorem 5, there exists a measurable  $f_i : [0, 1] \times \mathcal{X} \times \mathcal{M}(\mathcal{X}) \rightarrow \mathcal{X}$  such that

$$(\mathbf{X}_n, \mathbf{Y}_{n \setminus i}, Y_i) \stackrel{\text{a.s.}}{=} (\mathbf{X}_n, \mathbf{Y}_{n \setminus i}, f_i(\eta_i, X_i, \mathbb{M}_{\mathbf{X}_n})) ,$$

for  $\eta_i \sim \text{Unif}[0, 1]$  and  $\eta_i \perp\!\!\!\perp (\mathbf{X}_n, \mathbf{Y}_{n \setminus i})$ . This is true for each  $i \in [n]$ , and  $\mathbb{S}_n$ -equivariance implies that  $(\mathbf{X}_n, \mathbf{Y}_{n \setminus i}, Y_i) \stackrel{\text{d}}{=} (\mathbf{X}_n, Y_{[n] \setminus j}, Y_j)$  for all  $i, j \in [n]$ . Thus it is possible to choose the same function  $f_i = f$  for all  $i$ . This yields (23).

The reverse direction is easy to verify, since the noise variables are i.i.d.,  $\mathbf{X}_n$  is exchangeable, and  $\mathbb{M}_{\mathbf{X}_n}$  is permutation-invariant. ■

*The impact of the conditional independence assumption  $Y_i \perp\!\!\!\perp_{\mathbf{X}_n} (\mathbf{Y}_n \setminus Y_i)$ .* In the deterministic case, the assumed conditional independence relationships among the outputs  $\mathbf{Y}_n$  are trivially satisfied, so that (23) (without outsourced noise) is the most general form for a permutation-equivariant function. However, in the stochastic case, the assumed conditional independence significantly simplifies the structure of the conditional distribution and the corresponding functional representation. While the assumed conditional independence is key in the simplicity of the representation (23), it may limit the expressiveness: there are permutation-equivariant conditional distributions which do not satisfy the conditional independence assumption (examples are given below). On the other hand, without conditional independence between the elements of  $\mathbf{Y}_n$ , it is possible to show that  $Y_i = f(\eta_i, X_i, \mathbb{M}_{(X_j, Y_j)_{j \in [n] \setminus i}})$ . Such dependence between elements of  $\mathbf{Y}_n$ , although more expressive than (23), induces cycles in the computation graph, similar to Restricted Boltzmann Machines (Smolensky, 1987) and other Exponential Family Harmoniums (Welling et al., 2005). Furthermore, they may be limited in practice by the computational requirements of approximate inference algorithms. Striking a balance between flexibility and tractability via some simplifying assumption seems desirable.

Two examples illustrate the existence of permutation-equivariant conditional distributions which do not satisfy the conditional independence assumption made in Theorem 12, and suggest another assumption. Both examples have a similar structure: there exists some random variable, say  $W$ , such that the conditional independence  $Y_i \perp\!\!\!\perp_{(W, \mathbf{X}_n)} (\mathbf{Y}_n \setminus Y_i)$  holds. Assuming the existence of such a  $W$  would lead to the representation  $Y_i = f(\eta_i, W, X_i, \mathbb{M}_{\mathbf{X}_n})$ , and potentially allow for more expressive models, as  $W$  could be included in the neural network architecture and learned.

For the first example, let  $\mathbf{Y}_n$  be given as in (23), but with a vector of finitely exchangeable (but not i.i.d.) noise  $\boldsymbol{\eta}_n \perp\!\!\!\perp \mathbf{X}_n$ . Then  $\mathbf{Y}_n$  would still be conditionally  $\mathbb{S}_n$ -equivariant, but it would not satisfy the conditional independence assumption  $Y_i \perp\!\!\!\perp_{\mathbf{X}_n} (\mathbf{Y}_n \setminus Y_i)$ . However, it would satisfy  $Y_i \perp\!\!\!\perp_{(\mathbf{X}_n, \boldsymbol{\eta}_n)} (\mathbf{Y}_n \setminus Y_i)$ , which by similar arguments as in the proof of Theorem 12, implies the existence of a representation

$$Y_i = f'(\eta'_i, \eta_i, X_i, \mathbb{M}_{(X_i, \eta_i)_{i \in [n]}}), \quad (26)$$

for some other function  $f' : [0, 1]^2 \times \mathcal{X} \times \mathcal{M}(\mathcal{X} \times [0, 1])$  and i.i.d. noise  $\eta'_i \perp\!\!\!\perp (\mathbf{X}_n, \boldsymbol{\eta}_n)$ , in which case (23) would be a special case.

As a second example, in practice it is possible to construct more elaborate conditionally  $\mathbb{S}_n$ -equivariant distributions by composing multiple ones as in Theorem 2(ii). Suppose  $\mathbf{Y}_n$  is conditionally  $\mathbb{S}_n$ -equivariant and mutually independent given  $\mathbf{X}_n$ , and  $\mathbf{Z}_n$  is conditionally  $\mathbb{S}_n$ -equivariant and mutually independent given  $\mathbf{Y}_n$ . Theorem 2 shows that with  $\mathbf{Y}_n$  marginalized out,  $\mathbf{Z}_n$  is conditionally  $\mathbb{S}_n$ -equivariant given  $\mathbf{X}_n$ , while Theorem 12 guarantees the existence of the functional representations for each  $i \in [n]$ :

$$(Y_i)_{i \in [n]} = (f(\eta_i, X_i, \mathbb{M}_{\mathbf{X}_n}))_{i \in [n]} \quad \text{and} \quad (Z_i)_{i \in [n]} = (f'(\eta'_i, Y_i, \mathbb{M}_{\mathbf{Y}_n}))_{i \in [n]}.$$

Substituting the functional representation for  $\mathbf{Y}_n$  into that for  $\mathbf{Z}_n$  yields, for each  $i \in [n]$ ,

$$Z_i = f'(\eta'_i, f(\eta_i, X_i, \mathbb{M}_{\mathbf{X}_n}), \mathbb{M}_{(f(\eta_i, X_i, \mathbb{M}_{\mathbf{X}_n}))_{i \in [n]}}) = f''(\eta'_i, \eta_i, X_i, \mathbb{M}_{(X_i, \eta_i)_{i \in [n]}}),$$

for some  $f''$ , which is a special case of (26) and which implies  $Z_i \perp\!\!\!\perp_{(\mathbf{X}_n, \boldsymbol{\eta}_n)} (\mathbf{Z}_n \setminus Z_i)$ .

*Modeling  $\mathbb{S}_n$ -equivariance with neural networks.* One choice for  $\mathbb{S}_n$ -equivariant neural networks is

$$Y_i = f(\eta_i, X_i, g(\mathbf{X}_n)) \quad (27)$$

where  $f$  is an arbitrary (stochastic) neural network module, and  $g$  an arbitrary permutation-invariant module (say one of the examples in Section 6.1). An example of a permutation-equivariant module using a permutation-invariant submodule is shown in the middle panel of Fig. 4.

Theorem 2 enables the use of an architectural algebra for constructing complex permutation-invariant and -equivariant stochastic neural network modules. Specifically, permutation-equivariant modules may be constructed by composing simpler permutation-equivariant modules (see the middle panel of Fig. 4), while permutation-invariant modules can be constructed by composing permutation-equivariant modules with a final permutation-invariant module (right panel of Fig. 4). Some examples from the literature illustrate.

**Example 9 ( $\mathbb{S}_n$ -equivariant neural network layers, revisited)** *It is straightforward to see that Example 2 is a deterministic special case of Theorem 12:*

$$Y_i = \sigma(\theta_0 X_i + \theta_1 \sum_{j=1}^n X_j)$$

where  $\sum_{j=1}^n X_j = \int_{\mathcal{X}} \mathbb{M}_{\mathbf{X}_n}(dx)$  is a function of the empirical measure. While this example's nonlinear activation of linear combinations encodes a typical feed-forward neural network structure, Theorem 12 shows that more general functional relationships are allowed, for example (27).

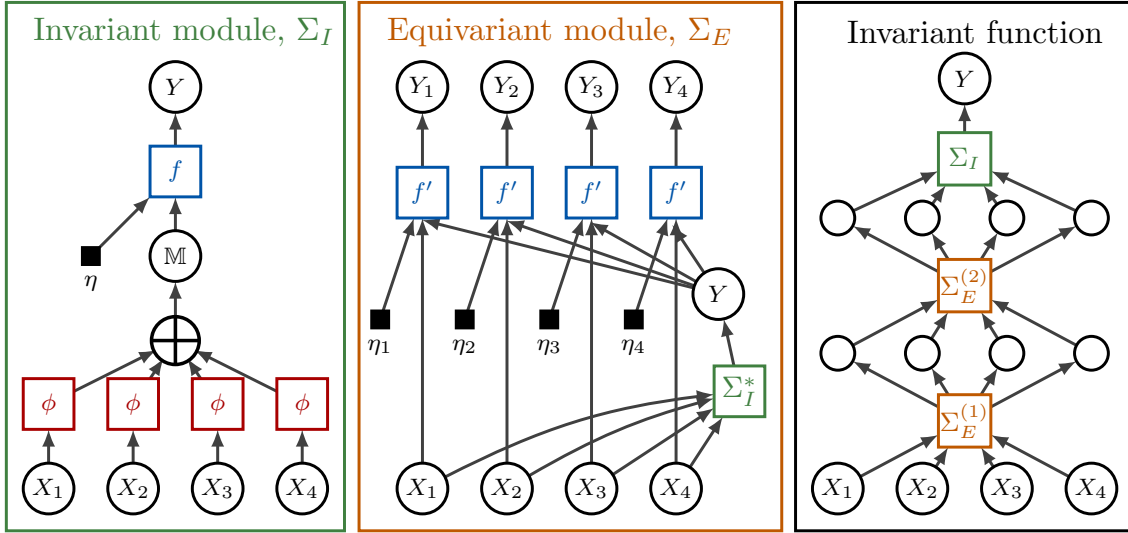


Figure 4: *Left*: An invariant module depicting (21). *Middle*: An equivariant module depicting (27); note that the invariant sub-module,  $\Sigma_I^*$ , must be deterministic unless there are alternative conditional independence assumptions, such as (26). *Right*: An invariant stochastic function composed of equivariant modules. Functional representations of  $\mathbb{S}_n$ -invariant and -equivariant conditional distributions. Circles denote random variables, with a row denoting an exchangeable sequence. The blue squares denote arbitrary functions, possibly with outsourced noise  $\eta$  which are mutually independent and independent of everything else. Same labels mean that the functions are the same. Red squares denote arbitrary embedding functions, possibly parameterized by a neural network, and  $\oplus$  denotes a symmetric pooling operation. Orange rectangles denote a module which gives a functional representation of a  $\mathbb{S}_n$ -equivariant conditional distribution. Likewise green rectangles for permutation-invariant conditional distributions.

**Example 10 (Equivariance and convolution)** *Kondor and Trivedi (2018) characterized the properties of deterministic feed-forward neural networks that are equivariant under the action of a compact group,  $\mathcal{G}$ . Roughly speaking, their results show that each layer  $\ell$  of the network must be a convolution of the output of the previous layer with some filter  $\chi_\ell$ . The general form of the convolution is defined in group theoretic terms that are beyond the scope of this paper. However, in the case of an exchangeable sequence, the situation is particularly simple. As an alternative to (23),  $Y_i$  may be represented by a function  $f'(\eta_i, X_i, \mathbb{M}_{\mathbf{X}_{n \setminus i}})$  (see (25) in the proof of Theorem 12), which makes clear the structure of the relationship between  $\mathbf{X}_n$  and  $\mathbf{Y}_n$ : element  $Y_i$  has a “receptive field” that focuses on  $X_i$  and treats the elements  $\mathbf{X}_{n \setminus i}$  as a structureless background field via  $\mathbb{M}_{\mathbf{X}_{n \setminus i}}$ . The dependence on the latter is invariant under permutations  $\pi' \in \mathbb{S}_{n-1}$ ; in group theoretic language,  $\mathbb{S}_{n-1}$  stabilizes  $i$  in  $[n]$ . That is, all permutations that fix  $i$  form an equivalence class. As such, for each  $i$  the index set  $[n]$  is in one-to-one correspondence with the set of equivalent permutations*

that either move the  $i$ th element to some other element (there are  $n - 1$  of these), or fix  $i$ . This is the quotient space  $\mathbb{S}_n/\mathbb{S}_{n-1}$ ; by the main result of Kondor and Trivedi (2018), any  $\mathbb{S}_n$ -equivariant feed-forward network with hidden layers all of size  $n$  must be composed of connections between layers  $X \mapsto Y$  defined by the convolution

$$Y_i = \sigma((X * \chi)_i) = \sigma\left(\sum_{j=1}^n X_{(i+j) \bmod n} \chi(j)\right), \quad \chi(j) = \delta_n(j)(\theta_0 + \theta_1) + \sum_{k=1}^{n-1} \delta_k(j)\theta_1.$$

The representation may be interpreted as a convolution of the previous layer’s output with a filter “centered” on element  $X_i$ , and is equivalent to that of Example 9.

**Example 11 (Self-attention)** Lee et al. (2019) proposed a  $\mathbb{S}_n$ -invariant architecture based on self-attention (Vaswani et al., 2017). For an input set of  $n$   $d$ -dimensional observations, the so-called Set Transformer combines attention over the  $d$  input dimensions with nonlinear functions of pairwise interactions between the input observations. In the simplest implementation, with no attention components, the Set Transformer computes in each network layer a nonlinear activation of the Gramian matrix  $\mathbf{X}_n \mathbf{X}_n^T$ ; the full architecture with attention is somewhat complicated, and we refer the reader to Lee et al. (2019). Furthermore, to combat prohibitive computational cost, a method inspired by inducing point techniques from the Gaussian Process literature (Snelson and Ghahramani, 2006) was introduced. In a range of experiments focusing on tasks that benefit from modeling dependence between elements of the set, such as clustering, the Set Transformers architecture out-performed architectures that did not include pairwise or higher-order interactions, like Deep Sets (Example 1).

### 6.3. Input Sets of Variable Size

In applications, the data set may consist of finitely exchangeable sequences of varying length. Theorems 11 and 12 are statements about input sequences of fixed length  $n$ . In general, they suggest that a separate function  $f_n$  needs to be learned for each  $n$  for which there is a corresponding observation  $\mathbf{X}_n$  in the data set. As a practical matter, clearly this is undesirable. In practice, the most common approach is to compute an independent embedding  $\phi : \mathcal{X} \rightarrow \mathbb{R}$  of each element of an input set, and then combine the embeddings with a symmetric pooling function like sum or max. For example,  $f(\mathbf{X}_n) = \max\{\phi(X_1), \dots, \phi(X_n)\}$ , or (6) from Examples 1 and 2. Clearly, such a function is invariant under permutations. However, recent work has explored the use of pairwise and higher-order interactions in the pooling function (the work by Murphy et al., 2019, mentioned in Example 6 is an example); empirical evidence indicates that the increased functional complexity results in higher model capacity and better performance for tasks that rely on modeling the dependence between elements in the input set.

Intuitively, more complex functions, such as those composed of pairwise (or higher-order) functions of an input sequence, give rise to higher-capacity models able to model more complicated forms of dependence between elements of an input sequence. In the context of exchangeability, this can be made more precise, as a difference between finitely and infinitely exchangeable sequences. In particular, let  $\mathbf{X}_n$  be the length  $n$  prefix of an infinitely exchangeable sequence  $\mathbf{X}_\infty$ . If a sequence of sufficient statistics  $S_n : \mathcal{X}^n \rightarrow \mathcal{S}_n$  exists, the conditionally i.i.d. representation (1) of the distribution of  $\mathbf{X}_\infty$  requires that they have the following properties (Freedman, 1962; Lauritzen, 1984, 1988):

(i) *symmetry under permutation*:

$$S_n(\pi \cdot \mathbf{X}_n) = S_n(\mathbf{X}_n) \quad \text{for all } \pi \in \mathbb{S}_n, \quad n \in \mathbb{N}; \quad (28)$$

(ii) *recursive computability*: for all  $n, m \in \mathbb{N}$  there are functions  $\psi_{n,m} : \mathcal{S}_n \times \mathcal{S}_m \rightarrow \mathcal{S}_{n+m}$  such that

$$S_{n+m}(\mathbf{X}_{n+m}) = \psi_{n,m}(S_n(\mathbf{X}_n), S_m(\mathbf{X}_{n+m} \setminus \mathbf{X}_n)). \quad (29)$$

A statistic that satisfies these properties must be of the form (Lauritzen, 1988)

$$S_n(\mathbf{X}_n) = S_1(X_1) \oplus \cdots \oplus S_1(X_n), \quad (30)$$

where  $(\mathcal{S}_1, \oplus)$  is an Abelian group or semigroup. Equivalently, because of symmetry under permutation, we write  $S_n(\mathbb{M}_{\mathbf{X}_n})$ . Examples with  $\mathcal{X} = \mathbb{R}_+$  include:

- (i)  $S_1(X_i) = \log X_i$  with  $S_1(X_i) \oplus S_1(X_j) = \log X_i + \log X_j$ ;
- (ii)  $S_1(X_i) = X_i$  with  $S_1(X_i) \oplus S_1(X_j) = X_i \vee X_j$ ;
- (iii)  $S_1(X_i) = \delta_{X_i}(\bullet)$  with  $S_1(X_i) \oplus S_1(X_j) = \delta_{X_i}(\bullet) + \delta_{X_j}(\bullet)$ .

Observe that pairwise (i.e., second-order) and higher-order statistics that do not decompose into first-order functions are precluded by the recursive computability property; an example is  $S_n(\mathbb{M}_{\mathbf{X}_n}) = \sum_{i,j \in [n]} X_i X_j$ .

Infinite exchangeability restricts the dependence between elements in  $\mathbf{X}_n$ : using a model based on only first-order functions, so that the properties of permutation symmetry and recursive computability are satisfied, limits the types of dependence that the model can capture. In practice, this can lead to shortcomings. For example, an infinitely exchangeable sequence cannot have negative correlation  $\rho = \text{Corr}(X_i, X_j)$ , but a finitely exchangeable sequence that cannot be extended to a longer exchangeable sequence can have  $\rho < 0$  (e.g., Aldous, 1985, pp. 7-8). One way to interpret this fact is that the type of dependence that gives rise to negative covariance cannot be captured by first-order functions. When using  $\mathbf{X}_n$  to predict another random variable, the situation becomes more complex, but to the extent that adequate (sufficient and d-separating) statistics are used, the same concepts are relevant. Ultimately, a balance between flexibility and computational efficiency must be found; the exact point of balance will depend on the details of the problem and may require novel computational methods (i.e., the inducing point-like methods in Example 11, or the sampling approximations in Murphy et al. (2019)) so that more flexible function classes can be used.

#### 6.4. Partially Exchangeable Sequences and Layers of Different Sizes

Shawe-Taylor (1989) and Ravanbakhsh et al. (2017) consider the problem of input-output invariance under a general discrete group  $\mathcal{G}$  acting on a standard feed-forward network, which consists of layers, potentially with different numbers of nodes, connected by weights. Those papers each found that  $\mathcal{G}$  and the neural network architecture must form a compatible pair: a pair of layers, treated as a weighted bipartite graph, forms a  $\mathcal{G}$ -equivariant function



if the action of  $\mathcal{G}$  on that graph is an automorphism. Essentially,  $\mathcal{G}$  must partition the nodes of each layer into weight-preserving orbits; Ravanbakhsh et al. (2017) provide some illuminating examples.

In the language of exchangeability, such a partition of elements corresponds to *partial exchangeability*:<sup>9</sup> the distributional invariance of  $\mathbf{X}_{n_x}$  under the action of a subgroup of the symmetric group,  $\mathcal{G} \subset \mathbb{S}_n$ . Partially exchangeable analogues of Theorems 11 and 12 are possible using the same types of conditional independence arguments. The resulting functional representations would express elements of  $\mathbf{Y}_{n_y}$  in terms of the empirical measures of blocks in the partition of  $\mathbf{X}_{n_x}$ ; the basic structure is already present in Theorem 9, but the details would depend on conditional independence assumptions among the elements of  $\mathbf{Y}_{n_y}$ . We omit a specific statement of the result, which would require substantial notational development, for brevity.

## 7. Learning From Finitely Exchangeable Matrices and Graphs

Neural networks that operate on graph-valued input data have been useful for a range of tasks, from molecular design (Duvenaud et al., 2015) and quantum chemistry (Gilmer et al., 2017), to knowledge-base completion (Hamaguchi et al., 2017).

In this section, we consider random matrices<sup>10</sup> whose distribution is invariant to permutations applied to the index set. In particular, let  $\mathbf{X}_{\mathbf{n}_2}$  be a two-dimensional  $\mathcal{X}$ -valued array with index set  $[\mathbf{n}_2] := [n_1] \times [n_2]$ , such that  $X_{i,j}$  is the element of  $\mathbf{X}_{\mathbf{n}_2}$  at position  $(i, j)$ . Let  $\pi_k \in \mathbb{S}_{n_k}$  be a permutation of the set  $[n_k]$  for  $k \in \{1, 2\}$ . Denote by  $\mathbb{S}_{\mathbf{n}_2}$  the direct product  $\mathbb{S}_{n_1} \times \mathbb{S}_{n_2}$ . A collection of permutations  $\boldsymbol{\pi}_2 := (\pi_1, \pi_2) \in \mathbb{S}_{\mathbf{n}_2}$  acts on  $\mathbf{X}_{\mathbf{n}_2}$  in the natural way, separately on the corresponding dimension:

$$[\boldsymbol{\pi}_2 \cdot \mathbf{X}_{\mathbf{n}_2}]_{i,j} = X_{\pi_1(i), \pi_2(j)}. \quad (31)$$

The distribution of  $\mathbf{X}_{\mathbf{n}_2}$  is *separately exchangeable* if

$$\boldsymbol{\pi}_2 \cdot \mathbf{X}_{\mathbf{n}_2} = (X_{\pi_1(i), \pi_2(j)})_{i \in [n_1], j \in [n_2]} \stackrel{d}{=} (X_{i,j})_{i \in [n_1], j \in [n_2]} = \mathbf{X}_{\mathbf{n}_2}, \quad (32)$$

for every collection of permutations  $\boldsymbol{\pi}_2 \in \mathbb{S}_{\mathbf{n}_2}$ . We say that  $\mathbf{X}_{\mathbf{n}_2}$  is separately exchangeable if its distribution is.

For symmetric arrays, such that  $n_1 = n_2 = n$  and  $X_{i,j} = X_{j,i}$ , a different notion of exchangeability is needed. The distribution of a symmetric  $\mathcal{X}$ -valued array  $\bar{\mathbf{X}}_{\mathbf{n}}$  is *jointly exchangeable* if, for all  $\pi \in \mathbb{S}_n$ ,

$$\pi \cdot \bar{\mathbf{X}}_{\mathbf{n}} = (\bar{X}_{\pi(i), \pi(j)})_{i, j \in [n]} \stackrel{d}{=} (\bar{X}_{i,j})_{i, j \in [n]} = \bar{\mathbf{X}}_{\mathbf{n}}. \quad (33)$$

### 7.1. Sufficient Representations of Exchangeable Matrices

In order to obtain functional representation results for matrices, a suitable analogue to the empirical measure is required. In contrast to the completely unstructured empirical

9. Diaconis and Freedman (1984) use partial exchangeability to mean any number of probabilistic symmetries; we use it only to mean partial permutation-invariance.

10. The results here are special cases of general results for  $d$ -dimensional arrays. For simplicity, we present the two-dimensional case and consider the general case in Appendix C.

measure of a sequence defined in (11), a sufficient representation of a matrix must retain the structural information encoded by the rows and columns of the matrix, but discard any ordering information. For matrices, such an object corresponds to a step function, also called the empirical graphon or a checkerboard function (Orbanz and Roy, 2015; Borgs and Chayes, 2017), which has played a central role in the theory of large exchangeable graphs and their representations as graphons (Lovász, 2012).

In the context of the current work, working with the checkerboard function is unnecessarily complicated; any maximal invariant of the symmetric group acting on  $\mathcal{X}$ -valued arrays will suffice. Specifically, any *canonical form*  $\mathbf{C}_{\mathbf{X}_{n_2}} := \text{CANON}(\mathbf{X}_{n_2})$  (or  $\mathbf{C}_{\overline{\mathbf{X}}_n} = \text{CANON}(\overline{\mathbf{X}}_n)$  for symmetric matrices), which serves as an orbit representative, can be used. The computational problem of finding a canonical representative of a graph is known as *canonical labeling*, and is at least as hard as the graph isomorphism problem (which is in NP, though it is unknown whether it is in P or NP-complete). In practice, some graph automorphism tools rely on canonical labeling algorithms, of which NAUTY and TRACES (McKay and Piperno, 2014) have demonstrated remarkably good performance. Canonical labeling results in a permutation  $\pi_2$  applied to the input to obtain an orbit representative, and therefore the process also can be used to define a representative equivariant  $\tau : \mathcal{X}^{n_1 \times n_2} \rightarrow \mathbb{S}_{n_2}$ . Different input domains  $\mathcal{X}$  may admit different canonicalization procedures, though for most practical purposes,  $\mathcal{X} \subseteq \mathbb{R}$ .

## 7.2. $\mathbb{S}_{n_2}$ -Invariant Conditional Distributions

As in Section 6.1, consider modeling  $Y$  as the output of a function whose input is a separately exchangeable matrix  $\mathbf{X}_{n_2}$ . In analogy to sequences in Theorem 11, sufficiency of  $\mathbf{C}_{\mathbf{X}_{n_2}}$  characterizes the class of all finitely exchangeable distributions on  $\mathcal{X}^{n_1 \times n_2}$ . The proof simply requires showing that  $\mathbf{C}_{\mathbf{X}_{n_2}}$  is a maximal invariant of  $\mathbb{S}_{n_2}$  acting on  $\mathcal{X}^{n_1 \times n_2}$ . To state the result, for any fixed canonicalization procedure, let  $\mathcal{C}_{n_2}(\mathcal{X})$  denote the space of canonical forms of matrices in  $\mathcal{X}^{n_1 \times n_2}$ .

**Theorem 13** *Let  $\mathbf{X}_{n_2}$  be a separately exchangeable  $\mathcal{X}$ -valued matrix indexed by  $[n_1] \times [n_2]$ , and  $Y \in \mathcal{Y}$  another random variable. Then  $P_{Y|\mathbf{X}_{n_2}}$  is  $\mathbb{S}_{n_2}$ -invariant if and only if there is a measurable function  $f : [0, 1] \times \mathcal{C}_{n_2}(\mathcal{X}) \rightarrow \mathcal{Y}$  such that*

$$(\mathbf{X}_{n_2}, Y) \stackrel{\text{a.s.}}{=} (\mathbf{X}_{n_2}, f(\eta, \mathbf{C}_{\mathbf{X}_{n_2}})) \quad \text{where } \eta \sim \text{Unif}[0, 1] \quad \text{and } \eta \perp\!\!\!\perp \mathbf{X}_{n_2}. \quad (34)$$

Furthermore,  $\mathbf{C}_{\mathbf{X}_{n_2}}$  is sufficient for the family  $\mathcal{P}_{\mathbf{X}_{n_2}}^{\mathbb{S}_{n_2}}$  and adequate for  $\mathcal{P}_{\mathbf{X}_{n_2}, Y}^{\mathbb{S}_{n_2}}$ .

**Proof**  $\mathbf{C}_{\mathbf{X}_{n_2}}$  is a maximal invariant by construction, and therefore Theorems 7 and 10 yield the result. ■

An identical result holds for jointly exchangeable matrices  $\overline{\mathbf{X}}_n$ , with symmetric canonical form  $\mathbf{C}_{\overline{\mathbf{X}}_n}$ .

As was the case for sequences, any function of a canonical form is invariant under permutations of the index set of  $\mathbf{X}_{n_2}$ . Most of the recent examples from the deep learning literature incorporate vertex features; that composite case is addressed in Section 7.4. A simple example without vertex features is the read-out layer of a neural network that operates on undirected, symmetric graphs.

**Example 12 (Read-out for message-passing neural networks)** *Gilmer et al. (2017)* reviewed recent work on neural networks whose input is an undirected graph (i.e., a symmetric matrix) on vertex set  $[n]$ , and whose hidden layers act as message-passing operations between vertices of the graph. These message-passing hidden layers are equivariant (see Sections 7.3 and 7.4); adding a final invariant layer makes the whole network invariant. A particularly simple architecture involves a single message-passing layer and no input features on the vertices, and a typical permutation-invariant read-out layer of the form

$$R = \sum_{i \in [n]} f \left( \sum_{j \in [n]} h(X_{i,j}) \right) = f'(\mathbf{C}_{\bar{\mathbf{X}}_n}).$$

### 7.3. $\mathbb{S}_{\mathbf{n}_2}$ -Equivariant Conditional Distributions

In analogy to  $\mathbf{X}_n$  and  $\mathbf{Y}_n$  in Section 6.2,  $\mathbf{X}_{\mathbf{n}_2}$  and  $\mathbf{Y}_{\mathbf{n}_2}$  might represent adjacent neural network layers; in such cases the goal is to transfer the symmetry of  $\mathbf{X}_{\mathbf{n}_2}$  to  $\mathbf{Y}_{\mathbf{n}_2}$ , and permutation-equivariance is the property of interest. With a collection of permutations acting on  $\mathbf{X}_{\mathbf{n}_2}$  and  $\mathbf{Y}_{\mathbf{n}_2}$  as in (31), permutation-equivariance is defined in the same way as for sequences. In particular, if  $\mathbf{X}_{\mathbf{n}_2}$  is exchangeable then  $\mathbf{Y}_{\mathbf{n}_2}$  is conditionally  $\mathbb{S}_{\mathbf{n}_2}$ -equivariant if and only if

$$(\pi_2 \cdot \mathbf{X}_{\mathbf{n}_2}, \pi_2 \cdot \mathbf{Y}_{\mathbf{n}_2}) \stackrel{d}{=} (\mathbf{X}_{\mathbf{n}_2}, \mathbf{Y}_{\mathbf{n}_2}) \quad \text{for all } \pi_2 \in \mathbb{S}_{\mathbf{n}_2}. \quad (35)$$

The main result in this section is a functional representation of  $\mathbf{Y}_{\mathbf{n}_2}$  in terms of a separately exchangeable array  $\mathbf{X}_{\mathbf{n}_2}$ , when the elements of  $\mathbf{Y}_{\mathbf{n}_2}$  are also conditionally independent given  $\mathbf{X}_{\mathbf{n}_2}$ .<sup>11</sup> In particular, each element  $Y_{i,j}$  is expressed in terms of  $X_{i,j}$ , outsourced noise  $\eta_{i,j}$ , and an augmented canonical form defined as follows.

Let  $\mathbf{C}_{i,:}$  denote the  $i$ th row of  $\mathbf{C}_{\mathbf{X}_{\mathbf{n}_2}}$ , and  $\mathbf{C}_{:,j}$  the  $j$ th column; define the *separately augmented canonical form* as

$$[\mathbf{C}_{\mathbf{X}_{\mathbf{n}_2}}^{(i,j)}]_{k,\ell} = ([\mathbf{C}_{\mathbf{X}_{\mathbf{n}_2}]_{k,\ell}, [\mathbf{C}_{\mathbf{X}_{\mathbf{n}_2}]_{i,\ell}, [\mathbf{C}_{\mathbf{X}_{\mathbf{n}_2}]_{k,j}]). \quad (36)$$

$\mathbf{C}_{\mathbf{X}_{\mathbf{n}_2}}^{(i,j)}$  augments the canonical form  $\mathbf{C}_{\mathbf{X}_{\mathbf{n}_2}}$  with the  $i$ th row and  $j$ th column, which are broadcast over the appropriate dimensions; one encodes  $\mathbf{C}_{i,:}$  and one encodes  $\mathbf{C}_{:,j}$ . These are analogues of the empirical measures  $\mathbb{M}_{\mathbf{C}_{i,:}}$  and  $\mathbb{M}_{\mathbf{C}_{:,j}}$ , but with their structure coupled to that of  $\mathbf{C}_{\mathbf{X}_{\mathbf{n}_2}}$ . Denote by  $\mathcal{C}_{\mathbf{n}_2}^{\text{aug}}(\mathcal{X})$  the space of all such functions augmented canonical forms of matrices with dimension  $\mathbf{n}_2$ .

A general version of the following theorem, for  $d$ -dimensional arrays, is given in Appendix C. The proof has the same basic structure as the one for sequences in Theorem 12, but with substantially more notation. Below, the proof for  $d = 2$  is given in order to highlight the important structure.

**Theorem 14** *Suppose  $\mathbf{X}_{\mathbf{n}_2}$  and  $\mathbf{Y}_{\mathbf{n}_2}$  are  $\mathcal{X}$ - and  $\mathcal{Y}$ -valued arrays, respectively, each indexed by  $[n_1] \times [n_2]$ , and that  $\mathbf{X}_{\mathbf{n}_2}$  is separately exchangeable. Assume that the elements of  $\mathbf{Y}_{\mathbf{n}_2}$  are mutually conditionally independent given  $\mathbf{X}_{\mathbf{n}_2}$ . Then  $\mathbf{Y}_{\mathbf{n}_2}$  is conditionally  $\mathbb{S}_{\mathbf{n}_2}$ -equivariant*

11. This is satisfied by neural networks without intra-layer or skip connections. Weaker conditional independence assumptions may be considered, as in Section 6.2.

given  $\mathbf{X}_{\mathbf{n}_2}$  if and only if there is a measurable function  $f : [0, 1] \times \mathcal{X} \times \mathcal{C}_{\mathbf{n}_2}^{\text{aug}}(\mathcal{X}) \rightarrow \mathcal{Y}$  such that

$$\left( \mathbf{X}_{\mathbf{n}_2}, \mathbf{Y}_{\mathbf{n}_2} \right) \stackrel{\text{a.s.}}{=} \left( \mathbf{X}_{\mathbf{n}_2}, \left( f(\eta_{i,j}, X_{i,j}, \mathbf{C}_{\mathbf{X}_{\mathbf{n}_2}}^{(i,j)}) \right)_{i \in [n_1], j \in [n_2]} \right), \quad (37)$$

for i.i.d. uniform random variables  $(\eta_{i,j})_{i \in [n_1], j \in [n_2]} \perp\!\!\!\perp \mathbf{X}_{\mathbf{n}_2}$ .

**Proof** First, assume that  $\mathbf{Y}_{\mathbf{n}_2}$  is conditionally  $\mathbb{S}_{\mathbf{n}_2}$ -equivariant given  $\mathbf{X}_{\mathbf{n}_2}$ . Then  $\pi_2 \cdot (\mathbf{X}_{\mathbf{n}_2}, \mathbf{Y}_{\mathbf{n}_2}) \stackrel{\text{d}}{=} (\mathbf{X}_{\mathbf{n}_2}, \mathbf{Y}_{\mathbf{n}_2})$  for all  $\pi_2 \in \mathbb{S}_{\mathbf{n}_2}$ . Let  $\mathbb{S}_{\mathbf{n}_2}^{(i,j)} \subset \mathbb{S}_{\mathbf{n}_2}$  be the stabilizer subgroup of  $(i, j)$ , i.e., the set of permutations that fixes element  $(i, j)$  in  $\mathbf{X}_{\mathbf{n}_2}$ . Note that each  $\pi_2^{(i,j)} \in \mathbb{S}_{\mathbf{n}_2}^{(i,j)}$  fixes both  $X_{i,j}$  and  $Y_{i,j}$ , and that  $\mathbb{S}_{\mathbf{n}_2}^{(i,j)}$  is homomorphic to  $\mathbb{S}_{\mathbf{n}_2-1}$ . Observe that any  $\pi_2^{(i,j)} \in \mathbb{S}_{\mathbf{n}_2}^{(i,j)}$  may rearrange the elements within the  $i$ th row of  $\mathbf{X}_{\mathbf{n}_2}$ , but it remains the  $i$ th row in  $\pi_2^{(i,j)} \cdot \mathbf{X}_{\mathbf{n}_2}$ . Similarly, the elements in the  $j$ th column,  $\mathbf{X}_{:,j}$  may be rearranged but remains the  $j$ th column. As a result, the  $j$ th element of every row is fixed (though it moves with its row), as is the  $i$ th element of every column.

That fixed-element structure will be used to establish the necessary conditional independence relationships. To that end, let  $r_i : [n_1] \setminus i \rightarrow [n_1 - 1]$  map the row indices of  $\mathbf{X}_{\mathbf{n}_2}$  to the row indices of the matrix obtained by removing the  $i$ th row from  $\mathbf{X}_{\mathbf{n}_2}$ :

$$r_i^{-1}(k) = \begin{cases} k & k < i \\ k - 1 & k > i \end{cases}.$$

Analogously, let  $c_j^{-1} : [n_2] \setminus j \rightarrow [n_2 - 1]$  map the column indices of  $\mathbf{X}_{\mathbf{n}_2}$  to those of  $\mathbf{X}_{\mathbf{n}_2}$  with the  $j$ th column removed. Define the  $\mathcal{X}^3$ -valued array  $\mathbf{Z}^{(i,j)}$  as

$$[\mathbf{Z}^{(i,j)}]_{k,\ell} = (X_{r_i(k), c_j(\ell)}, X_{i, c_j(\ell)}, X_{r_i(k), j}), \quad k \in [n_1 - 1], \ell \in [n_2 - 1]. \quad (38)$$

That is,  $\mathbf{Z}^{(i,j)}$  is formed by removing  $\mathbf{X}_{i,:}$  and  $\mathbf{X}_{:,j}$  from  $\mathbf{X}_{\mathbf{n}_2}$  (and  $X_{i,j}$  from  $\mathbf{X}_{i,:}$  and  $\mathbf{X}_{:,j}$ ), and broadcasting the removed row and column entries over the corresponding rows and columns of the matrix that remains.  $\mathbf{Z}^{(i,j)}$  inherits the exchangeability of  $\mathbf{X}_{\mathbf{n}_2}$  in the first element of each entry, and the fixed-elements structure in the second two elements, and therefore overall it is separately exchangeable:

$$\pi'_2 \cdot \mathbf{Z}^{(i,j)} \stackrel{\text{d}}{=} \mathbf{Z}^{(i,j)}, \quad \text{for all } \pi'_2 \in \mathbb{S}_{\mathbf{n}_2-1}.$$

Now, marginally (for  $Y_{i,j}$ ),

$$(\pi'_2 \cdot \mathbf{Z}^{(i,j)}, (X_{i,j}, Y_{i,j})) \stackrel{\text{d}}{=} (\mathbf{Z}^{(i,j)}, (X_{i,j}, Y_{i,j})), \quad \text{for all } \pi'_2 \in \mathbb{S}_{\mathbf{n}_2-1}.$$

Therefore, by Theorem 13,  $(X_{i,j}, Y_{i,j}) \perp\!\!\!\perp_{\mathbf{C}_{\mathbf{Z}^{(i,j)}}} \mathbf{Z}^{(i,j)}$ . Conditioning on  $X_{i,j}$  and  $\mathbf{C}_{\mathbf{Z}^{(i,j)}}$  is equivalent to conditioning on  $X_{i,j}$  and  $\mathbf{C}_{\mathbf{X}_{\mathbf{n}_2}}^{(i,j)}$ , yielding

$$Y_{i,j} \perp\!\!\!\perp_{(X_{i,j}, \mathbf{C}_{\mathbf{X}_{\mathbf{n}_2}}^{(i,j)})} \mathbf{X}_{\mathbf{n}_2}. \quad (39)$$

By Theorem 5, there is a measurable function  $f_{i,j} : [0, 1] \times \mathcal{X} \times \mathcal{C}_{\mathbf{n}_2}^{\text{aug}} \rightarrow \mathcal{Y}$  such that

$$Y_{i,j} = f_{i,j}(\eta_{i,j}, X_{i,j}, \mathbf{C}_{\mathbf{X}_{\mathbf{n}_2}}^{(i,j)}),$$

for a uniform random variable  $\eta_{i,j} \perp\!\!\!\perp \mathbf{X}_{\mathbf{n}_2}$ . This is true for all  $i \in [n_1]$  and  $j \in [n_2]$ ; by equivariance the same  $f_{i,j}$  must work for every  $(i, j)$ . Furthermore, by assumption the elements of  $\mathbf{Y}_{\mathbf{n}_2}$  are mutually conditionally independent given  $\mathbf{X}_{\mathbf{n}_2}$ , and therefore by the chain rule for conditional independence (Kallenberg, 2002, Prop. 6.8), the joint identity (37) holds.

The reverse direction is straightforward to verify. ■

A particularly simple version of (37) is

$$Y_{i,j} = f\left(\eta_{i,j}, X_{i,j}, \sum_{k=1}^{n_2} h_1(X_{i,k}), \sum_{\ell=1}^{n_1} h_2(X_{\ell,j}), \sum_{k,\ell} h_3(X_{\ell,k})\right), \quad (40)$$

for some functions  $h_m : \mathcal{X} \rightarrow \mathbb{R}$ ,  $m \in \{1, 2, 3\}$ . Clearly, this is conditionally equivariant. The following example from the deep learning literature is an even simpler version.

**Example 13 (Array-based MLPs)** *Hartford et al. (2018) determined the parameter-sharing schemes that result from deterministic permutation-equivariant MLP layers for matrices. They call such layers “exchangeable matrix layers”.<sup>12</sup> The equivariant weight-sharing scheme yields a simple expression:*

$$Y_{i,j} = \sigma\left(\theta_0 + \theta_1 X_{i,j} + \theta_2 \sum_{k=1}^{n_2} X_{i,k} + \theta_3 \sum_{\ell=1}^{n_1} X_{\ell,j} + \theta_4 \sum_{k,\ell} X_{\ell,k}\right).$$

That is,  $Y_{i,j}$  is a nonlinear activation of a linear combination of  $X_{i,j}$ , the sums of the  $j$ th column and  $i$ th row of  $\mathbf{X}_{\mathbf{n}_2}$ , and the sum of the entire matrix  $\mathbf{X}_{\mathbf{n}_2}$ . It is straightforward to see that this is a special deterministic case of (40). Hartford et al. (2018) also derive analogous weight-sharing schemes for MLPs for  $d$ -dimensional arrays; those correspond with the  $d$ -dimensional version of Theorem 14 (see Appendix C).

*Jointly exchangeable arrays.* The case of jointly exchangeable symmetric arrays is of particular interest because it applies to graph-valued data. Importantly, the edge variables are not restricted to be  $\{0, 1\}$ -valued; in practice they often take values in  $\mathbb{R}_+$ .

Let  $\overline{\mathbf{X}}_{\mathbf{n}}$  be a jointly exchangeable matrix in  $\mathcal{X}^{n \times n}$ , and  $\mathbf{C}_{\overline{\mathbf{X}}_{\mathbf{n}}}$  its canonical form. Similarly to (36), define the *jointly augmented canonical form* for a symmetric array as

$$[\mathbf{C}_{\overline{\mathbf{X}}_{\mathbf{n}}}^{\{i,j\}}]_{k,\ell} = ([\mathbf{C}_{\overline{\mathbf{X}}_{\mathbf{n}}}]_{k,\ell}, \{([\mathbf{C}_{\overline{\mathbf{X}}_{\mathbf{n}}}]_{i,k}, [\mathbf{C}_{\overline{\mathbf{X}}_{\mathbf{n}}}]_{j,k}), ([\mathbf{C}_{\overline{\mathbf{X}}_{\mathbf{n}}}]_{i,\ell}, [\mathbf{C}_{\overline{\mathbf{X}}_{\mathbf{n}}}]_{j,\ell})\}). \quad (41)$$

The symmetry of  $\overline{\mathbf{X}}_{\mathbf{n}}$  requires that the row and column entries be paired, which results in the second, set-valued, element on the right-hand side of (41). (The curly braces indicate a

<sup>12</sup>. This is a misnomer: exchangeability is a distributional property and there is nothing random.

set that is insensitive to the order of its elements, as opposed to parentheses, which indicate a sequence that is sensitive to order.) Denote by  $\overline{\mathcal{C}}_{\mathbf{n}}^{\text{aug}}(\mathcal{X})$  the space of all such jointly augmented canonical forms on  $\mathcal{X}^{n \times n}$ . The following counterpart of Theorem 14 applies to jointly exchangeable arrays such as undirected graphs.

**Theorem 15** *Suppose  $\overline{\mathbf{X}}_{\mathbf{n}}$  and  $\overline{\mathbf{Y}}_{\mathbf{n}}$  are symmetric  $\mathcal{X}$ - and  $\mathcal{Y}$ -valued arrays, respectively, each indexed by  $[n] \times [n]$ , and that  $\overline{\mathbf{X}}_{\mathbf{n}}$  is jointly exchangeable. Assume that the elements of  $\overline{\mathbf{Y}}_{\mathbf{n}}$  are mutually conditionally independent given  $\overline{\mathbf{X}}_{\mathbf{n}}$ . Then  $\overline{\mathbf{Y}}_{\mathbf{n}}$  is conditionally  $\mathbb{S}_n$ -equivariant given  $\overline{\mathbf{X}}_{\mathbf{n}}$  if and only if there is a measurable function  $f : [0, 1] \times \mathcal{X} \times \overline{\mathcal{C}}_{\mathbf{n}}^{\text{aug}}(\mathcal{X}) \rightarrow \mathcal{Y}$  such that*

$$(\overline{\mathbf{X}}_{\mathbf{n}}, \overline{\mathbf{Y}}_{\mathbf{n}}) \stackrel{\text{a.s.}}{=} \left( \overline{\mathbf{X}}_{\mathbf{n}}, (f(\eta_{i,j}, \overline{X}_{i,j}, \overline{\mathbf{C}}_{\overline{\mathbf{X}}_{\mathbf{n}}}^{\{i,j\}}))_{i \in [n], j \in [n]} \right), \quad (42)$$

for i.i.d. uniform random variables  $(\eta_{i,j})_{i \in [n], j \leq i} \perp\!\!\!\perp \overline{\mathbf{X}}_{\mathbf{n}}$  with  $\eta_{i,j} = \eta_{j,i}$ .

**Proof** The proof is essentially the same as that of Theorem 14. The main difference stems from the symmetry of  $\overline{\mathbf{X}}_{\mathbf{n}}$ : fixing row  $i$  also fixes column  $i$ , and fixing column  $j$  also fixes row  $j$ . Hence, the augmentations to the matrix consist of the paired sequences  $(\overline{\mathbf{X}}_{i,:}, \overline{\mathbf{X}}_{:,j}) = ((\overline{X}_{i,k}, \overline{X}_{j,k}))_{k \in [n]}$ . The counterpart to  $\mathbf{Z}^{(i,j)}$  in (38) is

$$[\overline{\mathbf{Z}}^{\{i,j\}}]_{k,\ell} = (\overline{X}_{k,\ell}, \{(\overline{X}_{i,\ell}, \overline{X}_{j,\ell}), (\overline{X}_{k,i}, \overline{X}_{k,j})\}), \quad k, \ell \in [n-1].$$

It is straightforward to show that  $\overline{\mathbf{Z}}^{\{i,j\}}$  is jointly exchangeable; the rest of the argument is the same as in the proof of Theorem 14.  $\blacksquare$

#### 7.4. Vertex Features: Equivariant Functions on the Vertex Set

In many applications of neural networks to graph-structured data, the desired output is a function defined on the vertex set, which can be thought of as a collection of vertex features. Additionally, the input graph may include vertex features. Formally, these features are encoded as a sequence  $\mathbf{X}_n$  whose elements correspond to the rows and columns of a symmetric matrix  $\overline{\mathbf{X}}_{\mathbf{n}}$ . We assume for simplicity that both  $\mathbf{X}_n$  and  $\overline{\mathbf{X}}_{\mathbf{n}}$  are  $\mathcal{X}$ -valued, but it is not strictly necessary that they take values in the same domain. In the separately exchangeable case, two distinct feature sequences  $\mathbf{X}_n^{(1)}$  and  $\mathbf{X}_n^{(2)}$  correspond to rows and columns, respectively. For simplicity, we focus on the symmetric case. A permutation  $\pi \in \mathbb{S}_n$  acts simultaneously on  $\mathbf{X}_n$  and  $\overline{\mathbf{X}}_{\mathbf{n}}$ , and the feature-augmented graph is jointly exchangeable if

$$(\pi \cdot \mathbf{X}_n, \pi \cdot \overline{\mathbf{X}}_{\mathbf{n}}) \stackrel{\text{d}}{=} (\mathbf{X}_n, \overline{\mathbf{X}}_{\mathbf{n}}), \quad \text{for all } \pi \in \mathbb{S}_n. \quad (43)$$

Likewise, we may define a vertex feature sequence  $\mathbf{Y}_n$ , whose elements correspond to the rows and columns of  $\overline{\mathbf{Y}}_{\mathbf{n}}$ .

Consider the array  $\mathbf{B}(\mathbf{X}_n, \overline{\mathbf{X}}_{\mathbf{n}})$  obtained by broadcasting the vertex features over the appropriate dimensions of  $\overline{\mathbf{X}}_{\mathbf{n}}$ :  $[\mathbf{B}(\mathbf{X}_n, \overline{\mathbf{X}}_{\mathbf{n}})]_{i,j} = (\overline{X}_{i,j}, \{X_i, X_j\})$ . In this case, we denote the canonical form by  $\mathbf{C}_{\mathbf{B}(\mathbf{X}_n, \overline{\mathbf{X}}_{\mathbf{n}})}$ . Observe that  $\mathbf{B}(\mathbf{X}_n, \overline{\mathbf{X}}_{\mathbf{n}})$  is a jointly exchangeable

array because  $(\mathbf{X}_n, \overline{\mathbf{X}}_n)$  is exchangeable as in (43). Furthermore, if  $\overline{\mathbf{Y}}_n$  is conditionally  $\mathbb{S}_n$ -equivariant given  $\overline{\mathbf{X}}_n$ , then  $\mathbf{B}(\mathbf{Y}_n, \overline{\mathbf{Y}}_n)$ , with entries  $[\mathbf{B}(\mathbf{Y}_n, \overline{\mathbf{Y}}_n)]_{i,j} = (\overline{Y}_{i,j}, \{Y_i, Y_j\})$ , is conditionally  $\mathbb{S}_n$ -equivariant given  $\mathbf{B}(\mathbf{X}_n, \overline{\mathbf{X}}_n)$ . However, Theorem 15 does not apply:  $\mathbf{B}(\mathbf{Y}_n, \overline{\mathbf{Y}}_n)$  does not satisfy the mutually conditionally independent elements assumption because, for example,  $Y_i$  appears in every element in the  $i$ th row.

Further assumptions are required to obtain a useful functional representation of  $\overline{\mathbf{Y}}_n$ . In addition to the mutual conditional independence of the elements of  $\overline{\mathbf{Y}}_n$  given  $(\mathbf{X}_n, \overline{\mathbf{X}}_n)$ , we assume further that

$$Y_i \perp\!\!\!\perp_{(\mathbf{X}_n, \overline{\mathbf{X}}_n, \overline{\mathbf{Y}}_n)} (\mathbf{Y}_n \setminus Y_i), \quad i \in [n]. \quad (44)$$

In words, the elements of the output vertex feature sequence are conditionally independent, given the input data  $\mathbf{X}_n, \overline{\mathbf{X}}_n$  and the output edge features  $\overline{\mathbf{Y}}_n$ . This is consistent with the implicit assumptions used in practice in the deep learning literature, and leads to a representation with simple structure.

To state the result, let  $\mathbf{C}_{\mathbf{B}(\mathbf{X}_n, \overline{\mathbf{X}}_n)}$  be the augmented canonical form that includes the vertex features,

$$[\mathbf{C}_{\mathbf{B}(\mathbf{X}_n, \overline{\mathbf{X}}_n)}^{\{i,j\}}]_{k,\ell} = (\overline{X}_{k,\ell}, \{X_k, X_\ell\}, \{(\overline{X}_{i,\ell}, \overline{X}_{j,\ell}), (\overline{X}_{k,i}, \overline{X}_{k,j})\}),$$

which belongs to the space denoted  $\mathcal{C}_{\mathbf{B},\mathbf{n}}^{\text{aug}}(\mathcal{X})$ . Let  $\mathbf{S}(\overline{\mathbf{X}}_n, \overline{\mathbf{Y}}_n)$  be the symmetric (stacked) array with entries  $(\overline{X}_{i,j}, \overline{Y}_{i,j})$ , and  $\mathbf{B}(\mathbf{X}_n, \mathbf{S}(\overline{\mathbf{X}}_n, \overline{\mathbf{Y}}_n))$  be the same, with the entries of  $\mathbf{X}_n$  broadcast.

**Theorem 16** *Suppose  $(\mathbf{X}_n, \overline{\mathbf{X}}_n)$  and  $(\mathbf{Y}_n, \overline{\mathbf{Y}}_n)$  are  $\mathcal{X}$ - and  $\mathcal{Y}$ -valued vertex feature-augmented arrays, and that  $(\mathbf{X}_n, \overline{\mathbf{X}}_n)$  is jointly exchangeable as in (43). Assume that the elements of  $\overline{\mathbf{Y}}_n$  are mutually conditionally independent given  $(\mathbf{X}_n, \overline{\mathbf{X}}_n)$ , and that  $\mathbf{Y}_n$  satisfies (44). Then  $(\mathbf{Y}_n, \overline{\mathbf{Y}}_n)$  is conditionally  $\mathbb{S}_n$ -equivariant given  $(\mathbf{X}_n, \overline{\mathbf{X}}_n)$  if and only if there are measurable functions  $f_e : [0, 1] \times \mathcal{X} \times \overline{\mathcal{X}}^2 \times \mathcal{C}_{\mathbf{B},\mathbf{n}}^{\text{aug}}(\mathcal{X}) \rightarrow \mathcal{Y}$  and  $f_v : [0, 1] \times \mathcal{X} \times \mathcal{X} \times \mathcal{Y} \times \mathcal{C}_{\mathbf{B},\mathbf{n}}^{\text{aug}}(\mathcal{X} \times \mathcal{Y}) \rightarrow \mathcal{Y}$  such that*

$$\overline{Y}_{i,j} \stackrel{\text{a.s.}}{=} f_e(\eta_{i,j}, \overline{X}_{i,j}, \{X_i, X_j\}, \mathbf{C}_{\mathbf{B}(\mathbf{X}_n, \overline{\mathbf{X}}_n)}^{\{i,j\}}), \quad i, j \in [n] \quad (45)$$

$$Y_i \stackrel{\text{a.s.}}{=} f_v(\eta_i, X_i, \overline{X}_{i,i}, \overline{Y}_{i,i}, \mathbf{C}_{\mathbf{B}(\mathbf{X}_n, \mathbf{S}(\overline{\mathbf{X}}_n, \overline{\mathbf{Y}}_n))}^{\{i\}}), \quad i \in [n], \quad (46)$$

for i.i.d. uniform random variables  $(\eta_{i,j})_{i \in [n], j \leq i} \perp\!\!\!\perp (\mathbf{X}_n, \overline{\mathbf{X}}_n)$  with  $\eta_{i,j} = \eta_{j,i}$ , and  $(\eta_i)_{i \in [n]} \perp\!\!\!\perp (\mathbf{X}_n, \overline{\mathbf{X}}_n, \overline{\mathbf{Y}}_n)$ .

**Proof** The proof, like that of Theorem 15, is essentially the same as for Theorem 14. Incorporating vertex features requires that for any permutation  $\pi \in \mathbb{S}_n$ , the fixed elements of  $\mathbf{X}_n$  be collected along with the fixed rows and columns of  $\overline{\mathbf{X}}_n$ ; the structure of the argument is identical.  $\blacksquare$

Equations (45) and (46) indicate that given an input  $(\mathbf{X}_n, \overline{\mathbf{X}}_n)$  and functional forms for  $f_e$  and  $f_v$ , computation of  $\overline{\mathbf{Y}}_n$  and  $\mathbf{Y}_n$  proceeds in two steps: first, compute the elements  $\overline{Y}_{i,j}$  of  $\overline{\mathbf{Y}}_n$ ; second, compute the vertex features  $\mathbf{Y}_n$  from  $\mathbf{X}_n, \overline{\mathbf{X}}_n$ , and  $\overline{\mathbf{Y}}_n$ . Note that

within each step, computations can be parallelized due to the conditional independence assumptions.

The following examples from the literature are special cases of Theorem 16.

**Example 14 (Graph-based structured prediction)** *Herzig et al. (2018)* considered the problem of deterministic permutation-equivariant structured prediction in the context of mapping images to scene graphs. In particular, for a weighted graph with edge features  $\bar{\mathbf{X}}_n$  and vertex features  $\mathbf{X}_n$ , those authors define a graph labeling function  $f$  to be “graph-permutation invariant” (GPI) if  $f(\pi \cdot (\mathbf{X}_n, \bar{\mathbf{X}}_n)) = \pi \cdot f(\mathbf{X}_n, \bar{\mathbf{X}}_n)$  for all  $\pi \in \mathbb{S}_n$ .<sup>13</sup> Furthermore, they implicitly set  $\bar{\mathbf{Y}}_n = \bar{\mathbf{X}}_n$ , and assume that  $\tilde{X}_{i,i} = 0$  for all  $i \in [n]$  and that  $\mathbf{X}_n$  is included in  $\bar{\mathbf{X}}_n$  (e.g., on the diagonal). The main theoretical result is that a graph labeling function is GPI if and only if

$$[f(\mathbf{X}_n, \bar{\mathbf{X}}_n)]_i = \rho(X_i, \sum_j \alpha(X_j, \sum_k \phi(X_j, \tilde{X}_{j,k}, X_k))) ,$$

for appropriately defined functions  $\rho$ ,  $\alpha$ , and  $\phi$ . Inspection of the proof reveals that the second argument of  $\rho$  (the sum of  $\alpha$  functions) is equivalent to a maximal invariant. In experiments, a particular GPI neural network architecture showed better sample efficiency for training, as compared with an LSTM with the inputs in random order and with a fully connected feed-forward network, each network with the same number of parameters.

**Example 15 (Message passing graph neural networks)** *Gilmer et al. (2017)* reviewed the recent literature on graph-based neural network architectures, and found that many of them fit in the framework of so-called message passing neural networks (MPNNs). MPNNs take as input a graph with vertex features and edge features (the features may be vector-valued, but for simplicity of notation we assume they are scalar-valued real numbers). Each neural network layer  $\ell$  acts as a round of message passing between adjacent vertices, with the typically edge-features held fixed from layer to layer. In particular, denote the (fixed) input edge features by  $\bar{\mathbf{X}}_n$  and the computed vertex features of layer  $\ell - 1$  by  $\mathbf{X}_n$ ; then the vertex features for layer  $\ell$ ,  $\mathbf{Y}_n$ , are computed as

$$Y_i = U_\ell(X_i, \sum_{j \in [n]} \mathbb{1}_{\{\tilde{X}_{i,j} > 0\}} M_\ell(X_i, X_j, \tilde{X}_{i,j})) , \quad i \in [n] ,$$

for some functions  $M_\ell : \mathbb{R}^3 \rightarrow \mathbb{R}$  and  $U_\ell : \mathbb{R}^2 \rightarrow \mathbb{R}$ . This is a deterministic special case of (46). *Gilmer et al. (2017)* note that *Kearnes et al. (2016)* also compute different edge features,  $\tilde{Y}_{i,j}$ , in each layer. The updated edge features are used in place of  $\tilde{X}_{i,j}$  in the third argument of  $M_\ell$  in the equation above, and are an example of the two-step implementation of (45)-(46).

*Separate exchangeability with features.* Separately exchangeable matrices with features, which may be regarded as representing a bipartite graph with vertex features, have a similar representation. Let  $\mathbf{X}_n^{(1)}$  and  $\mathbf{X}_n^{(2)}$  denote the row and column features, respectively. Separate exchangeability for such data structures is defined in the obvious way. Furthermore, a similar conditional independence assumption is made for the output vertex features:

13. Although *Herzig et al. (2018)* call this invariance, we note that it is actually equivariance.



$\mathbf{Y}_n^{(1)} \perp\!\!\!\perp_{(\mathbf{X}_n^{(1)}, \mathbf{X}_n^{(2)}, \mathbf{X}_{n_2}, \mathbf{Y}_{n_2})} \mathbf{Y}_n^{(2)}$  and

$$Y_i^{(m)} \perp\!\!\!\perp_{(\mathbf{X}_n^{(1)}, \mathbf{X}_n^{(2)}, \mathbf{X}_{n_2}, \mathbf{Y}_{n_2})} (\mathbf{Y}_n^{(m)} \setminus Y_i^{(m)}), \quad i \in [n_m], m \in \{1, 2\}. \quad (47)$$

The representation result is stated here for completeness. To state the result, note that appropriate broadcasting of features requires that  $\mathbf{X}_n^{(1)}$  be broadcast over columns  $\mathbf{X}_{:,j}$ , and  $\mathbf{X}_n^{(2)}$  be broadcast over rows  $\mathbf{X}_{i,:}$ . We denote this with the same broadcasting operator as above,  $\mathbf{B}(\mathbf{X}_n^{(1)}, \mathbf{X}_n^{(2)}, \mathbf{X}_{n_2})$ . The proof is similar to that of Theorem 16, and is omitted for brevity.

**Theorem 17** *Suppose  $(\mathbf{X}_n^{(1)}, \mathbf{X}_n^{(2)}, \mathbf{X}_{n_2})$  and  $(\mathbf{Y}_n^{(1)}, \mathbf{Y}_n^{(2)}, \mathbf{Y}_{n_2})$  are  $\mathcal{X}$ - and  $\mathcal{Y}$ -valued vertex feature-augmented arrays, and that  $(\mathbf{X}_n^{(1)}, \mathbf{X}_n^{(2)}, \mathbf{X}_{n_2})$  is separately exchangeable. Assume that the elements of  $\mathbf{Y}_{n_2}$  are mutually conditionally independent given  $(\mathbf{X}_n^{(1)}, \mathbf{X}_n^{(2)}, \mathbf{X}_{n_2})$ , and that  $\mathbf{Y}_n^{(1)}$  and  $\mathbf{Y}_n^{(2)}$  satisfy (47). Then  $(\mathbf{Y}_n^{(1)}, \mathbf{Y}_n^{(2)}, \mathbf{Y}_{n_2})$  is conditionally  $\mathbb{S}_{n_2}$ -equivariant given  $(\mathbf{X}_n^{(1)}, \mathbf{X}_n^{(2)}, \mathbf{X}_{n_2})$  if and only if there are measurable functions  $f_e : [0, 1] \times \mathcal{X}^3 \times \mathcal{C}_{\mathbf{B}, n_2}^{\text{aug}}(\mathcal{X}) \rightarrow \mathcal{Y}$  and  $f_v^{(m)} : [0, 1] \times \mathcal{X} \times \mathcal{X}^{n_m} \times \mathcal{Y}^{n_m} \times \mathcal{C}_{\mathbf{B}, n_2}^{\text{aug}}(\mathcal{X} \times \mathcal{Y}) \rightarrow \mathcal{Y}$ , for  $m \in \{1, 2\}$ , such that*

$$Y_{i,j} \stackrel{\text{a.s.}}{=} f_e(\eta_{i,j}, X_{i,j}, X_i^{(1)}, X_j^{(2)}, \mathbf{C}_{\mathbf{B}}(\mathbf{s}(\mathbf{X}_n^{(1)}, \mathbf{X}_{:,j}), \mathbf{s}(\mathbf{X}_n^{(2)}, \mathbf{X}_{i,:}), \mathbf{X}_{n_2})), \quad i \in [n_1], j \in [n_2] \quad (48)$$

$$Y_i^{(1)} \stackrel{\text{a.s.}}{=} f_v^{(1)}(\eta_i^{(1)}, X_i^{(1)}, \mathbf{X}_{i,:}, \mathbf{Y}_{i,:}, \mathbf{C}_{\mathbf{B}}(\mathbf{X}_n^{(1)}, \mathbf{s}(\mathbf{X}_n^{(2)}, \mathbf{X}_{i,:}), \mathbf{s}(\mathbf{X}_{n_2}, \mathbf{Y}_{n_2}))), \quad i \in [n_1], \quad (49)$$

$$Y_j^{(2)} \stackrel{\text{a.s.}}{=} f_v^{(2)}(\eta_j^{(2)}, X_j^{(2)}, \mathbf{X}_{:,j}, \mathbf{Y}_{:,j}, \mathbf{C}_{\mathbf{B}}(\mathbf{s}(\mathbf{X}_n^{(1)}, \mathbf{X}_{:,j}), \mathbf{X}_n^{(2)}, \mathbf{s}(\mathbf{X}_{n_2}, \mathbf{Y}_{n_2}))), \quad j \in [n_2], \quad (50)$$

for i.i.d. uniform random variables  $((\eta_{i,j})_{i \in [n_1], j \in [n_2]}) \perp\!\!\!\perp (\mathbf{X}_n^{(1)}, \mathbf{X}_n^{(2)}, \mathbf{X}_{n_2})$  and  $((\eta_i^{(1)})_{i \in [n_1]}, (\eta_j^{(2)})_{j \in [n_2]}) \perp\!\!\!\perp (\mathbf{X}_n^{(1)}, \mathbf{X}_n^{(2)}, \mathbf{X}_{n_2}, \mathbf{Y}_{n_2})$ .

## 8. Discussion

The probabilistic approach to symmetry has allowed us to draw on tools from an area that is typically outside the purview of the deep learning community. Those tools shed light on the underlying structure of previous work on invariant neural networks, and expand the scope of what is possible with such networks. Moreover, those tools place invariant neural networks in a broader statistical context, making connections to the fundamental concepts of sufficiency and adequacy.

To conclude, we give some examples of other data structures to which the theory developed in the present work could be applied and describe some questions prompted by this work and by others.

### 8.1. Other Exchangeable Structures

The results in Section 6 can be adapted to other exchangeable structures in a straightforward manner. We briefly describe two settings; Orbanz and Roy (2015) survey a number of other exchangeable structures from statistics and machine learning to which this work could apply.

*Edge-exchangeable graphs.* Cai et al. (2016); Williamson (2016); Crane and Dempsey (2018) specified generative models for network data as an exchangeable sequence of edges,  $\mathbf{E}_n = ((u, v)_1, \dots, (u, v)_n)$ . The sequence theory from Section 6 applies, rather than the array theory from Section 7, which applies to vertex-exchangeable graphs. However, incorporating vertex features into edge-exchangeable models would require some extra work, as a permutation acting on the edge sequence has a different (random) induced action on the vertices; we leave it as an interesting problem for future work.

The model of Caron and Fox (2017) is finitely edge-exchangeable when conditioned on the random number of edges in the network. Therefore, the sequence theory could be incorporated into inference procedures for that model, with the caveat that the neural network architecture would be required to accept inputs of variable size and the discussion from Section 6.3 would be relevant.

*Markov chains and recurrent processes.* A Markov chain  $\mathbf{Z}_n := (Z_1, Z_2, \dots, Z_n)$  on state-space  $\mathcal{Z}$  has exchangeable sub-structure. In particular, define a *z-block* as a sub-sequence of  $\mathbf{Z}_n$  that starts and ends on some state  $z \in \mathcal{Z}$ . Clearly, the joint probability

$$P_{\mathbf{Z}_n}(\mathbf{Z}_n) = P_{Z_1}(Z_1) \prod_{i=2}^n Q(Z_i | Z_{i-1})$$

is invariant under all permutations of  $Z_1$ -blocks (and of any other  $z$ -blocks for  $z \in \mathbf{Z}_n$ ). Denote these blocks by  $(B_{Z_1, j})_{j \geq 1}$ . Diaconis and Freedman (1980b) used this notion of *Markov exchangeability* to show that all recurrent processes on countable state-spaces are mixtures of Markov chains, and Bacallado et al. (2013) used similar ideas to analyze reversible Markov chains on uncountable state-spaces.

For a finite Markov chain, the initial state,  $Z_1$ , and the  $Z_1$ -blocks are sufficient statistics. Equivalently,  $Z_1$  and the empirical measure of the  $m \leq n$   $Z_1$ -blocks,

$$\mathbb{M}_{\mathbf{Z}_n}(\bullet) = \sum_{j=1}^m \delta_{B_{Z_1, j}}(\bullet),$$

plays the same role as  $\mathbb{M}_{\mathbf{X}_n}$  for an exchangeable sequence. (If  $\mathbf{Z}_n$  is the prefix of an infinite, recurrent process, then  $Z_1$  and the empirical measure of transitions are sufficient.) It is clear that the theory from Section 6 can be adapted to accommodate Markov chains; indeed, Wqvist et al. (2019) used this idea to learn summary statistics from a Markov chain for use in Approximate Bayesian Computation.

## 8.2. Open Questions

*More flexible conditional independence assumptions.* The simplicity of results on functional representations of equivariant conditional distributions relied on conditional independence assumptions like  $Y_i \perp\!\!\!\perp_{\mathbf{X}_n} (\mathbf{Y}_n \setminus Y_i)$ . As discussed in Section 6.2, additional flexibility could be obtained by assuming the existence of a random variable  $W$  such that  $Y_i \perp\!\!\!\perp_{(\mathbf{X}_n, W)} (\mathbf{Y}_n \setminus Y_i)$  holds, and learning  $W$  for each layer.  $W$  could be interpreted as providing additional shared “context” for the input-output relationship, a construct that has been useful in recent work (Edwards and Storkey, 2017; Kim et al., 2019).

*Choice of pooling functions.* The results here are quite general, and say nothing about what other properties the function classes under consideration should have. For example, different symmetric pooling functions may lead to very different performance. This seems to be well understood in practice; for example, the pooling operation of element-wise mean of the top decile in Chan et al. (2018) is somewhat unconventional, but appears to have been chosen based on performance. Recent theoretical work by Xu et al. (2019) studies different pooling operations in the context of graph discrimination tasks with graph neural networks, a problem also considered from a slightly different perspective by Shawe-Taylor (1993). Wagstaff et al. (2019) take up the problem for functions defined on sets.

*Learning and generalization.* Our results leave open questions pertaining to learning and generalization. However, they point to some potential approaches, one of which we sketch here. Shawe-Taylor (1991, 1995) applied PAC theory to derive generalization bounds for feed-forward networks that employ a symmetry-induced weight-sharing scheme; these bounds are improvements on those for standard fully-connected multi-layer perceptrons. Weight-sharing might be viewed as a form of pre-training compression; recent work uses PAC-Bayes theory to demonstrate the benefits of compressibility of trained networks (Arora et al., 2018; Zhou et al., 2019), and Lyle et al. (2020) make some initial progress along these lines.

## Acknowledgments

The authors are grateful to the editor and three anonymous referees for detailed comments that helped improve the paper. The authors also thank Bruno Ribeiro for pointing out an error, and to Hongseok Yang for pointing out typos and an unnecessary condition for Theorem 9, in an earlier draft. BBR’s and YWT’s research leading to these results has received funding from the European Research Council under the European Union’s Seventh Framework Programme (FP7/2007-2013) ERC grant agreement no. 617071.

## Appendix A. Proof of Theorem 5

d-separation is a statement of conditional independence and therefore the proof of Theorem 5 is a straightforward application of the following basic result about conditional independence, which appears (with different notation) as Proposition 6.13 in Kallenberg (2002).

**Lemma 18 (Conditional independence and randomization)** *Let  $X, Y, Z$  be random elements in some measurable spaces  $\mathcal{X}, \mathcal{Y}, \mathcal{Z}$ , respectively, where  $\mathcal{Y}$  is Borel. Then  $Y \perp\!\!\!\perp_Z X$  if and only if  $Y \stackrel{\text{a.s.}}{=} f(\eta, Z)$  for some measurable function  $f : [0, 1] \times \mathcal{Z} \rightarrow \mathcal{Y}$  and some uniform random variable  $\eta \perp\!\!\!\perp (X, Z)$ .*

## Appendix B. Orbit Laws, Maximal Invariants, and Representative Equivariants: Proofs for Section 4

The results in Section 4 (Theorems 7, 9 and 10) rely on the disintegration of  $\mathcal{G}$ -invariant distributions  $P_X$  into a distribution over orbits and a distribution on each orbit generated by the normalized Haar measure (the orbit law). This disintegration, established in Theorem 20, is the mathematically precise version of the intuitive explanation given in Section 4, just before Theorem 7. We apply the result, along with a key conditional independence property that it implies, in order to obtain Theorems 7, 9 and 10.

We require the following definition, taken from Eaton (1989).

**Definition 19** *A measurable cross-section is a set  $\mathcal{C} \subset \mathcal{X}$  with the following properties:*

- (i)  $\mathcal{C}$  is measurable.
- (ii) For each  $x \in \mathcal{X}$ ,  $\mathcal{C} \cap \mathcal{G} \cdot x$  consists of exactly one point, say  $\mathcal{C}_x$ .
- (iii) The function  $t : \mathcal{X} \rightarrow \mathcal{C}$  defined by  $t(x) = \mathcal{C}_x$  is  $\mathcal{B}_{\mathcal{X}}$ -measurable when  $\mathcal{C}$  has  $\sigma$ -algebra  $\mathcal{B}_{\mathcal{C}} = \{B \cap \mathcal{C} \mid B \in \mathcal{B}_{\mathcal{X}}\}$ .

One may think of a measurable cross-section as consisting of a single representative from each orbit of  $\mathcal{X}$  under  $\mathcal{G}$ . For a compact group acting measurably on a Borel space, there always exists a (not necessarily unique) measurable cross-section. For non-compact groups or more general spaces, the existence of a measurable cross-section as defined above is not guaranteed, and more powerful technical tools are required (e.g. Andersson, 1982; Schindler, 2003; Kallenberg, 2017).

*Orbit laws.* Intuitively, conditioned on  $X$  being on a particular orbit, it should be possible to sample a realization of  $X$  by first sampling a random group element  $G$ , and then applying  $G$  to a representative element of the orbit. More precisely, let  $\lambda_{\mathcal{G}}$  be the normalized Haar measure of  $\mathcal{G}$ , which is the unique left- and right-invariant measure on  $\mathcal{G}$  such that  $\lambda_{\mathcal{G}}(\mathcal{G}) = 1$ . Let  $M : \mathcal{X} \rightarrow \mathcal{S}$  be a maximal invariant. For any  $m \in \text{range}(M)$ , let  $c_m = M^{-1}(m)$  be the element of  $\mathcal{C}$  for which  $M(c_m) = m$ ; such an element always exists. Note that the set of elements in  $\mathcal{S}$  that do not correspond to an orbit of  $\mathcal{X}$ ,  $\mathcal{S}_{\emptyset} := \{m; M^{-1}(m) = \emptyset\}$ , has measure zero under  $P_X \circ M^{-1}$ .

For any  $B \in \mathcal{B}_{\mathcal{X}}$  and  $m \in \mathcal{S} \setminus \mathcal{S}_{\emptyset}$ , define the *orbit law* as

$$\mathbb{U}_m^{\mathcal{G}}(B) = \int_{\mathcal{G}} \delta_{g \cdot c_m}(B) \lambda_{\mathcal{G}}(dg) = \lambda_{\mathcal{G}}(\{g; g \cdot c_m \in B\}). \quad (51)$$

(For  $m \in \mathcal{S}_\emptyset$ ,  $\mathbb{U}_m^\mathcal{G}(B) = 0$ .) Observe that, in agreement with the intuition above, a sample from the orbit law can be generated by sampling a random group element  $G \sim \lambda_\mathcal{G}$ , and applying it to  $c_m$ . The orbit law inherits the invariance of  $\lambda_\mathcal{G}$  and acts like the uniform distribution on the elements of the orbit, up to fixed points.<sup>14</sup> For any integrable function  $\varphi : \mathcal{X} \rightarrow \mathbb{R}$ , the expectation with respect to the orbit law is

$$\mathbb{U}_m^\mathcal{G}[\varphi] = \int_{\mathcal{X}} \varphi(x) \mathbb{U}_m^\mathcal{G}(dx) = \int_{\mathcal{G}} \varphi(g \cdot c_m) \lambda_\mathcal{G}(dg). \quad (52)$$

The orbit law arises in the disintegration of any  $\mathcal{G}$ -invariant distribution  $P_X$  as the fixed distribution on each orbit.

**Lemma 20** *Let  $\mathcal{X}$  and  $\mathcal{S}$  be Borel spaces,  $\mathcal{G}$  a compact group acting measurably on  $\mathcal{X}$ , and  $M : \mathcal{X} \rightarrow \mathcal{S}$  a maximal invariant on  $\mathcal{X}$  under  $\mathcal{G}$ . If  $X$  is a random element of  $\mathcal{X}$ , then its distribution  $P_X$  is  $\mathcal{G}$ -invariant if and only if*

$$P_X(X \in \bullet \mid M(X) = m) = \mathbb{U}_m^\mathcal{G}(\bullet) = q(\bullet, m), \quad (53)$$

for some Markov kernel  $q : \mathcal{B}_\mathcal{X} \times \mathcal{S} \rightarrow \mathbb{R}_+$ . If  $P_X$  is  $\mathcal{G}$ -invariant and  $Y$  is any other random variable, then  $P_{Y|X}$  is  $\mathcal{G}$ -invariant if and only if  $Y \perp\!\!\!\perp_{M(X)} X$ .

**Proof** For the first claim, a similar result appears in Eaton (1989, Ch. 4). For completeness, a slightly different proof is given here. Suppose that  $P_X$  is  $\mathcal{G}$ -invariant. Then  $(g \cdot X, M) \stackrel{d}{=} (X, M)$  for all  $g \in \mathcal{G}$ . By Fubini's theorem, the statement is true even for a random group element  $G \perp\!\!\!\perp X$ . Let  $G$  be a random element of  $\mathcal{G}$ , sampled from  $\lambda_\mathcal{G}$ . Then for any measurable function  $\varphi : \mathcal{X} \rightarrow \mathbb{R}_+$ , and for any set  $A \in \sigma(M)$ ,

$$\mathbb{E}[\varphi(X); A] = \mathbb{E}[\varphi(G \cdot X); A] = \mathbb{E}[\mathbb{E}[\varphi(G \cdot X) \mid X]; A] = \mathbb{E}[\mathbb{U}_{M(X)}^\mathcal{G}[\varphi]; A],$$

which establishes the first equality of (53); the second equality follows from Kallenberg (2002, Thm. 6.3) because  $\mathcal{S}$  is a Borel space, as is  $\mathcal{X}$ . Conversely, suppose that (53) is true.  $P_X(X \in \bullet) = \mathbb{E}[\mathbb{U}_M^\mathcal{G}(\bullet)]$ , where the expectation is taken with respect to the distribution of  $M$ .  $\mathbb{U}_M^\mathcal{G}(\bullet)$  is  $\mathcal{G}$ -invariant for all  $M \in \mathcal{S}$ , and therefore so is the marginal distribution  $P_X$ .

For the second claim, suppose  $Y$  is such that  $P_{Y|X}$  is  $\mathcal{G}$ -invariant, which by Theorem 1 is equivalent to  $(g \cdot X, Y) \stackrel{d}{=} (X, Y)$  for all  $g \in \mathcal{G}$ . The latter is equivalent to  $g \cdot (X, Y) \stackrel{d}{=} (X, Y)$ , with  $g \cdot Y = Y$  almost surely, for all  $g \in \mathcal{G}$ . Therefore,  $\mathcal{G} \cdot Y = \{Y\}$  almost surely, and  $\widetilde{M}(X, Y) := (M(X), Y)$  is a maximal invariant of  $\mathcal{G}$  acting on  $\mathcal{X} \times \mathcal{Y}$ . Therefore,  $\mathbb{U}_{\widetilde{M}(X, Y)}^\mathcal{G}(A \times B) = \mathbb{U}_{M(X)}^\mathcal{G}(A) \delta_Y(B)$ , for  $A \in \mathcal{B}_\mathcal{X}$  and  $B \in \mathcal{B}_\mathcal{Y}$ , and

$$P_X(X \in \bullet \mid M(X) = m, Y) = \mathbb{U}_m^\mathcal{G}(\bullet) = P_X(X \in \bullet \mid M(X) = m),$$

which implies  $Y \perp\!\!\!\perp_{M(X)} X$ . The converse is straightforward to check. ■

14. The orbit law only coincides with the uniform distribution on the orbit if the action of  $\mathcal{G}$  on the orbit is *regular*, which corresponds to the action being transitive and *free* (or fixed-point free). Since by definition the action of  $\mathcal{G}$  is transitive on each orbit, the orbit law is equivalent to the uniform distribution on the orbit if and only if  $\mathcal{G}$  acts freely on each orbit; if the orbit has any fixed points (i.e.,  $g \cdot x = x$  for some  $x$  in the orbit and  $g \neq e$ ), then the fixed points will have higher probability mass under the orbit law.

Note that (53) is equivalent to the disintegration

$$P_X(X \in \bullet) = \int_{\mathcal{S}} \mathbb{U}_m^{\mathcal{G}}(\bullet) \nu(dm) = \int_{\mathcal{C}} \mathbb{U}_{M(x)}^{\mathcal{G}}(\bullet) \mu(dx),$$

for some maximal invariant  $M : \mathcal{X} \rightarrow \mathcal{S}$  and probability measures  $\nu$  on  $\mathcal{S}$  and  $\mu$  on the measurable cross-section  $\mathcal{C}$ . This type of statement is more commonly encountered (e.g., Eaton, 1989, Ch. 4-5) than (53), but the form of (53) emphasizes the role of the orbit law as the conditional distribution of  $X$  given  $M(X)$ , which is central to the development of ideas in Section 3.

*Example: orbit law of an exchangeable sequence.* The orbit law of an exchangeable sequence  $\mathbf{X}_n$  is also known as the, *urn law*. It is defined as a probability measure on  $\mathcal{X}^n$  corresponding to permuting the elements of  $\mathbb{M}_{\mathbf{X}_n}$  according to a uniformly sampled random permutation:

$$\mathbb{U}_{\mathbb{M}_{\mathbf{X}_n}}^{\mathbb{S}_n}(\bullet) = \frac{1}{n!} \sum_{\pi \in \mathbb{S}_n} \delta_{\pi \cdot \mathbf{X}_n}(\bullet). \quad (54)$$

The urn law is so called because it computes the probability of sampling any sequence without replacement from an urn with  $n$  balls, each labeled by an element of  $\mathbf{X}_n$ .<sup>15</sup> Diaconis and Freedman (1980a) and Kallenberg (2005, Proposition 1.8) prove results (in different forms) similar to Theorem 11 that rely on the urn law.

### B.1. Invariant Conditional Distributions: Proof of Theorem 7

**Theorem 7** *Let  $X$  and  $Y$  be random elements of Borel spaces  $\mathcal{X}$  and  $\mathcal{Y}$ , respectively, and  $\mathcal{G}$  a compact group acting measurably on  $\mathcal{X}$ . Assume that  $P_X$  is  $\mathcal{G}$ -invariant, and pick a maximal invariant  $M : \mathcal{X} \rightarrow \mathcal{S}$ , with  $\mathcal{S}$  another Borel space. Then  $P_{Y|X}$  is  $\mathcal{G}$ -invariant if and only if there exists a measurable function  $f : [0, 1] \times \mathcal{S} \rightarrow \mathcal{Y}$  such that*

$$(X, Y) \stackrel{\text{a.s.}}{=} (X, f(\eta, M(X))) \quad \text{with } \eta \sim \text{Unif}[0, 1] \text{ and } \eta \perp\!\!\!\perp X. \quad (14)$$

**Proof** With the conditional independence relationship  $Y \perp\!\!\!\perp_{M(X)} X$  established in Theorem 20, Theorem 7 follows from Theorem 5. ■

### B.2. Representative Equivariants and Proof of Theorem 9

Recall that a representative equivariant is an equivariant function  $\tau : \mathcal{X} \rightarrow \mathcal{G}$  satisfying

$$\tau(g \cdot x) = g \cdot \tau(x) \quad \text{for each } g \in \mathcal{G}.$$

---

15. The metaphor is that of an urn with  $n$  balls, each labeled by an element of  $\mathbb{M}_{\mathbf{X}_n}$ ; a sequence is constructed by repeatedly picking a ball uniformly at random from the urn, without replacement. Variations of such a scheme can be considered, for example sampling with replacement, or replacing a sampled ball with two of the same label. See Mahmoud (2008) for an overview of the extensive probability literature studying such processes, including the classical Pólya urn and its generalizations (e.g., Blackwell and MacQueen, 1973), which play important roles in Bayesian nonparametrics.

Recall also that  $\tau_x := \tau(x)$  and that  $\tau_x^{-1}$  is the inverse element in  $\mathcal{G}$  of  $\tau_x$ . The following intermediate result (restated from Section 4) establishes two useful properties of  $\tau$ , which are used in the proof of Theorem 9.

**Lemma 21** *For a group  $\mathcal{G}$  acting measurably on Borel spaces  $\mathcal{X}$  and  $\mathcal{Y}$ , a representative equivariant  $\tau : \mathcal{X} \rightarrow \mathcal{G}$ , as defined in (15), has the following properties:*

- (i) *The function  $M_\tau : \mathcal{X} \rightarrow \mathcal{X}$  defined by  $M_\tau(x) = \tau_x^{-1} \cdot x$  is a maximal invariant.*
- (ii) *For any mapping  $b : \mathcal{X} \rightarrow \mathcal{Y}$ , the function*

$$f(x) = \tau_x \cdot b(\tau_x^{-1} \cdot x), \quad x \in \mathcal{X}, \quad (16)$$

*is  $\mathcal{G}$ -equivariant:  $f(g \cdot x) = g \cdot f(x)$ ,  $g \in \mathcal{G}$ .*

**Proof** Property (i) is proved in Eaton (1989, Ch. 2), as follows. Observe that for all  $g \in \mathcal{G}$  and  $x \in \mathcal{X}$ ,

$$M_\tau(g \cdot x) = \tau_{g \cdot x}^{-1} \cdot (g \cdot x) = (g \cdot \tau_x)^{-1} \cdot (g \cdot x) = \tau_x^{-1} \cdot g^{-1} \cdot g \cdot x = \tau_x^{-1} \cdot x = M_\tau(x),$$

so  $M_\tau(x)$  is invariant. Now suppose  $M_\tau(x_1) = M_\tau(x_2)$  for some  $x_1, x_2$ , and define  $g = \tau_{x_1} \cdot \tau_{x_2}^{-1}$ . Then

$$\begin{aligned} \tau_{x_1}^{-1} \cdot x_1 &= \tau_{x_2}^{-1} \cdot x_2 \\ x_1 &= (\tau_{x_1} \cdot \tau_{x_2}^{-1}) \cdot x_2 = g \cdot x_2. \end{aligned}$$

Therefore,  $x_1$  and  $x_2$  are in the same orbit and  $M_\tau$  is a maximal invariant.

For (ii), consider an arbitrary mapping  $b : \mathcal{X} \rightarrow \mathcal{Y}$ , and define the function  $f$  by  $\tau_x \cdot b(\tau_x^{-1} \cdot x)$ , as in (16). Then

$$\begin{aligned} f(g \cdot x) &= \tau_{g \cdot x} \cdot b(\tau_{g \cdot x}^{-1} \cdot g \cdot x) = \tau_{g \cdot x} \cdot b(\tau_x^{-1} \cdot x) = \tau_{g \cdot x} \cdot \tau_x^{-1} \cdot f(x) \\ &= g \cdot \tau_x \cdot \tau_x^{-1} \cdot f(x) = g \cdot f(x), \end{aligned}$$

where the equivariance of  $\tau_x$  is used repeatedly. ■

We note that the equivariance requirement on  $\tau$  can be relaxed, at the price of an extra condition required for  $f$  defined in (ii) above to be equivariant. Specifically, any function  $\bar{\tau} : \mathcal{X} \rightarrow \mathcal{G}$  that satisfies  $\bar{\tau}_{g \cdot x}^{-1} \cdot g \cdot x = \bar{\tau}_x^{-1} \cdot x$  (i.e., part (i) above) can be used in the same way as  $\tau$  to construct an equivariant  $f$ , with the condition that  $\mathcal{G}_x \subseteq \mathcal{G}_{f(x)}$ . Kallenberg (2005), Lemma 7.10 proves such a result (and the existence of  $\bar{\tau}$ ) in the case of a discrete group.

**Theorem 9** *Let  $\mathcal{G}$  be a compact group acting measurably on Borel spaces  $\mathcal{X}$  and  $\mathcal{Y}$ , such that there exists a measurable representative equivariant  $\tau : \mathcal{X} \rightarrow \mathcal{G}$  satisfying (15). Suppose  $P_X$  is  $\mathcal{G}$ -invariant. Then  $P_{Y|X}$  is  $\mathcal{G}$ -equivariant if and only if there exists a measurable  $\mathcal{G}$ -equivariant function  $f : [0, 1] \times \mathcal{X} \rightarrow \mathcal{Y}$  satisfying (17) such that*

$$(X, Y) \stackrel{\text{a.s.}}{=} (X, f(\eta, X)) \quad \text{for each } g \in \mathcal{G}, \text{ with } \eta \sim \text{Unif}[0, 1] \text{ and } \eta \perp\!\!\!\perp X.$$

**Proof** The proof of sufficiency closely follows the proof of Lemma 7.11 in Kallenberg (2005). It relies on proving that  $\tau_X^{-1} \cdot Y \perp\!\!\!\perp_{M_\tau} X$ , and applying Theorems 5 and 8.

Assume that  $g \cdot (X, Y) \stackrel{d}{=} (X, Y)$  for all  $g \in \mathcal{G}$ . Let  $M_\tau(x) = \tau_x^{-1} \cdot x$ , and for any  $x \in \mathcal{X}$ , let  $M'_{\tau,x} : \mathcal{Y} \rightarrow \mathcal{Y}$  be defined by  $M'_{\tau,x}(y) = \tau_x^{-1} \cdot y$ , for  $y \in \mathcal{Y}$ . As shown in the proof of Theorem 8,  $\tau_x \cdot \tau_{g \cdot x}^{-1} \cdot g = \tau_x^{-1}$  for all  $x \in \mathcal{X}$ . Therefore, the random elements  $M_\tau(X) \in \mathcal{X}$  and  $M'_{\tau,X}(Y) \in \mathcal{Y}$  satisfy

$$(M_\tau(X), M'_{\tau,X}(Y)) = \tau_x^{-1} \cdot (X, Y) = \tau_{g \cdot X}^{-1} \cdot g \cdot (X, Y), \quad \text{a.s., } g \in \mathcal{G}. \quad (55)$$

Now let  $G \sim \lambda_{\mathcal{G}}$  be a random element of  $\mathcal{G}$  such that  $G \perp\!\!\!\perp (X, Y)$ . Observe that because  $G \sim \lambda_{\mathcal{G}}$ ,  $g \cdot G \stackrel{d}{=} G \cdot g \stackrel{d}{=} G$  for all  $g \in \mathcal{G}$ . Furthermore, by assumption  $g \cdot (X, Y) \stackrel{d}{=} (X, Y)$ . Therefore, using Fubini's theorem,

$$G \cdot (X, Y) \stackrel{d}{=} (X, Y), \quad \text{and} \quad (G \cdot \tau_X^{-1}, X, Y) \stackrel{d}{=} (G, X, Y). \quad (56)$$

Using (55) and (56),

$$\begin{aligned} (X, M_\tau(X), M'_{\tau,X}(Y)) &= (X, \tau_X^{-1} \cdot X, \tau_X^{-1} \cdot Y) \\ &\stackrel{d}{=} (G \cdot X, \tau_{G \cdot X}^{-1} \cdot G \cdot X, \tau_{G \cdot X}^{-1} \cdot G \cdot Y) \\ &= (G \cdot X, \tau_X^{-1} \cdot X, \tau_X^{-1} \cdot Y) \\ &\stackrel{d}{=} (G \cdot \tau_X^{-1} \cdot X, \tau_X^{-1} \cdot X, \tau_X^{-1} \cdot Y) \\ &= (G \cdot M_\tau(X), M_\tau(X), M'_{\tau,X}(Y)) \end{aligned} \quad (57)$$

That is, jointly with the orbit representative  $M_\tau(X)$  and the representative equivariant applied to  $Y$ ,  $M'_{\tau,X}(Y)$ , the distribution of  $X$  is the same as if applying a random group element  $G \sim \lambda_{\mathcal{G}}$  to the orbit representative.

Trivially,  $M_\tau(X) \perp\!\!\!\perp_{M_\tau(X)} M'_{\tau,X}(Y)$ . Therefore,  $G \cdot M_\tau(X) \perp\!\!\!\perp_{M_\tau(X)} M'_{\tau,X}(Y)$  is implied by  $G \perp\!\!\!\perp (X, Y)$ . This conditional independence is transferred by the distributional equality in (57) to

$$X \perp\!\!\!\perp_{M_\tau(X)} M'_{\tau,X}(Y).$$

Therefore, by Theorem 5, there exists some measurable  $b : [0, 1] \times \mathcal{X} \rightarrow \mathcal{Y}$  such that  $M'_{\tau,X}(Y) = \tau_X^{-1} \cdot Y \stackrel{\text{a.s.}}{=} b(\eta, M_\tau(X))$  for  $\eta \sim \text{Unif}[0, 1]$  and  $\eta \perp\!\!\!\perp X$ . Applying  $\tau_X$ , Theorem 8 implies that

$$Y \stackrel{\text{a.s.}}{=} \tau_X \cdot b(\eta, M_\tau(X)) = \tau_X \cdot b(\eta, \tau_X^{-1} \cdot X) =: f(\eta, X)$$

is  $\mathcal{G}$ -equivariant in the second argument. This establishes sufficiency.

Conversely, for necessity, using (17) and the assumption that  $g \cdot X \stackrel{d}{=} X$  for all  $g \in \mathcal{G}$ ,

$$g \cdot (X, Y) = (g \cdot X, f(\eta, g \cdot X)) \stackrel{d}{=} (X, f(\eta, X)) = (X, Y).$$

■



### B.3. Proof of Theorem 10

**Theorem 10** *Let  $\mathcal{G}$  be a compact group acting measurably on standard Borel spaces  $\mathcal{X}$  and  $\mathcal{Y}$ , and let  $\mathcal{S}$  be another Borel space. Then:*

- (i) *Any maximal invariant  $M : \mathcal{X} \rightarrow \mathcal{S}$  on  $\mathcal{X}$  under  $\mathcal{G}$  is an adequate statistic of  $X$  for  $Y$  with respect to any model  $\mathcal{P}_{X,Y} \subseteq \mathcal{P}_{X,Y}^{\text{inv}}$ . In particular, to each  $P_{X,Y} \in \mathcal{P}_{X,Y}^{\text{inv}}$  corresponds a measurable  $\mathcal{G}$ -invariant function  $f : [0, 1] \times \mathcal{S} \rightarrow \mathcal{Y}$  as in (14).*
- (ii)  *$M_\tau$  as in Theorem 8 is an adequate statistic of  $X$  for  $\tau_X^{-1} \cdot Y$  with respect to any model  $\mathcal{P}_{X,Y} \subseteq \mathcal{P}_{X,\tau_X^{-1} \cdot Y}$ , and to each  $P_{X,Y} \in \mathcal{P}_{X,\tau_X^{-1} \cdot Y}$  there corresponds a measurable function  $b : [0, 1] \times \mathcal{X} \rightarrow \mathcal{Y}$  as in (18) such that  $\tau_X \cdot b(\eta, \tau_X^{-1} \cdot X)$  is  $\mathcal{G}$ -equivariant.*

**Proof** (i) Let  $M : \mathcal{X} \rightarrow \mathcal{S}$  be any maximal invariant on  $\mathcal{X}$  under  $\mathcal{G}$ . From Theorem 20, the conditional distribution of  $X$  given  $M(X) = m$  is equal to the orbit law  $\mathbb{U}_m^{\mathcal{G}}$ .  $\mathcal{S}$  is assumed to be Borel, so there is a Markov kernel (e.g., Kallenberg, 2002, Thm. 6.3)  $q_M : \mathcal{B}_{\mathcal{X}} \times \mathcal{S} \rightarrow \mathbb{R}_+$  such that  $q_M(\bullet, m) = \mathbb{U}_m^{\mathcal{G}}(\bullet)$ , and therefore  $M$  is a sufficient statistic for  $\mathcal{P}_X^{\text{inv}}$ . Moreover, also by Theorem 20,  $M$   $d$ -separates  $X$  and  $Y$  under any  $P_{X,Y} \in \mathcal{P}_{X,Y}^{\text{inv}}$ , and therefore  $M$  is also an adequate statistic. Theorem 5 implies the identity (14).

(ii) The proof follows a similar argument as in (i). Theorem 8 established that  $M_\tau(X) = \tau_X^{-1} \cdot X$  is a representative equivariant, so by Theorem 20,  $M_\tau$  is sufficient for  $\mathcal{P}_X^{\text{inv}}$ , a subset of which is obtained by marginalizing over  $Y$  in each distribution belonging to  $\mathcal{P}_{X,\tau_X^{-1} \cdot Y}$ . The proof of Theorem 9 also established that  $\tau_X^{-1} \cdot Y \perp\!\!\!\perp_{M_\tau(X)} X$ , which, along with sufficiency for the marginal model, is the definition of adequacy.  $\blacksquare$

## Appendix C. Representations of Exchangeable $d$ -Dimensional Arrays

The functional representations of conditionally  $\mathbb{S}_{n_2}$ -invariant (Theorem 13) and -equivariant (Theorem 14) distributions in Section 7 are special cases with  $d = 2$  of more general results for  $d$ -dimensional arrays.

Some notation is needed in order to state the results. For a fixed  $d \in \mathbb{N}$ , let  $\mathbf{X}_{\mathbf{n}_d}$  be a  $d$ -dimensional  $\mathcal{X}$ -valued array (called a  $d$ -array) with index set  $[\mathbf{n}_d] := [n_1] \times \cdots \times [n_d]$ , and with  $X_{i_1, \dots, i_d}$  the element of  $\mathbf{X}_{\mathbf{n}_d}$  at position  $(i_1, \dots, i_d)$ . The size of each dimension is encoded in the vector of integers  $\mathbf{n}_d = (n_1, \dots, n_d)$ . In this section, we consider random  $d$ -arrays whose distribution is invariant to permutations applied to the index set. In particular let  $\pi_k \in \mathbb{S}_{n_k}$  be a permutation of the set  $[n_k]$ . Denote by  $\mathbb{S}_{\mathbf{n}_d} := \mathbb{S}_{n_1} \times \cdots \times \mathbb{S}_{n_d}$  the direct product of each group  $\mathbb{S}_{n_k}$ ,  $k \in [d]$ .

A collection of permutations  $\boldsymbol{\pi}_d := (\pi_1, \dots, \pi_d) \in \mathbb{S}_{\mathbf{n}_d}$  acts on  $\mathbf{X}_{\mathbf{n}_d}$  in the natural way, separately on the corresponding input dimension of  $\mathbf{X}_{\mathbf{n}_d}$ :

$$[\boldsymbol{\pi}_d \cdot \mathbf{X}_{\mathbf{n}_d}]_{i_1, \dots, i_d} = X_{\pi_1(i_1), \dots, \pi_d(i_d)} \cdot$$

The distribution of  $\mathbf{X}_{\mathbf{n}_d}$  is separately exchangeable if

$$(X_{i_1, \dots, i_d})_{(i_1, \dots, i_d) \in [n_1, \dots, d]} \stackrel{d}{=} (X_{\pi_1(i_1), \dots, \pi_d(i_d)})_{(i_1, \dots, i_d) \in [n_1, \dots, d]}, \quad (58)$$

for every collection of permutations  $\pi_d \in \mathbb{S}_{\mathbf{n}_d}$ . We say that  $\mathbf{X}_{\mathbf{n}_d}$  is separately exchangeable if its distribution is.

For a symmetric array  $\overline{\mathbf{X}}_{\mathbf{n}_d}$ , such that  $n_1 = n_2 = \dots = n_d = n$  and  $\overline{X}_{i_1, \dots, i_d} = \overline{X}_{i_{\rho(1)}, \dots, i_{\rho(d)}}$  for all  $\rho \in \mathbb{S}_d$ , the distribution of  $\overline{\mathbf{X}}_{\mathbf{n}_d}$  is *jointly exchangeable* if, for all  $\pi \in \mathbb{S}_n$

$$(\overline{X}_{i_1, \dots, i_d})_{(i_1, \dots, i_d) \in [n]^d} \stackrel{d}{=} (\overline{X}_{\pi(i_1), \dots, \pi(i_d)})_{(i_1, \dots, i_d) \in [n]^d}. \quad (59)$$

As in Section 7.1, a canonical form of  $\mathbf{X}_{\mathbf{n}_d}$  is required. For a particular canonicalization routine, denote the space of  $d$ -dimensional canonical forms by  $\mathcal{C}_{\mathbf{n}_d}(\mathcal{X})$ .

The first result concerns  $\mathbb{S}_{\mathbf{n}_d}$ -invariant conditional distributions. An equivalent result holds for jointly exchangeable  $d$ -arrays and  $\mathbb{S}_n$ -invariant conditional distributions, which we omit for brevity.

**Theorem 22** *Suppose  $\mathbf{X}_{\mathbf{n}_d}$  is a separately exchangeable  $\mathcal{X}$ -valued array on index set  $\mathbf{n}_d$ , and  $Y \in \mathcal{Y}$  is another random variable. Then  $P_{Y|\mathbf{X}_{\mathbf{n}_d}}$  is  $\mathbb{S}_{\mathbf{n}_d}$ -invariant if and only if there is a measurable function  $f : [0, 1] \times \mathcal{C}_{\mathbf{n}_d}(\mathcal{X}) \rightarrow \mathcal{Y}$  such that*

$$(\mathbf{X}_{\mathbf{n}_d}, Y) \stackrel{\text{a.s.}}{=} (\mathbf{X}_{\mathbf{n}_d}, f(\eta, \mathbf{C}_{\mathbf{X}_{\mathbf{n}_d}})) \quad \text{where } \eta \sim \text{Unif}[0, 1] \quad \text{and } \eta \perp\!\!\!\perp \mathbf{X}_{\mathbf{n}_d}. \quad (60)$$

**Proof** Clearly, the representation (60) satisfies  $(\pi_d \cdot \mathbf{X}_{\mathbf{n}_d}, Y) \stackrel{d}{=} (\mathbf{X}_{\mathbf{n}_d}, Y)$ , for all  $\pi_d \in \mathbb{S}_{\mathbf{n}_d}$ , which by Theorem 1 and the assumption that  $\mathbf{X}_{\mathbf{n}_d}$  is separately exchangeable implies that  $Y$  is conditionally  $\mathbb{S}_{\mathbf{n}_d}$ -invariant given  $\mathbf{X}_{\mathbf{n}_d}$ . The converse is a easy consequence of the fact that the canonical form is a maximal invariant of  $\mathbb{S}_{\mathbf{n}_d}$  acting on  $\mathcal{X}^{n_1 \times \dots \times n_d}$  and Theorem 10.  $\blacksquare$

*$\mathbb{S}_{\mathbf{n}_d}$ -equivariant conditional distributions.* Let  $\mathbf{Y}_{\mathbf{n}_d}$  be a  $\mathcal{Y}$ -valued array indexed by  $[\mathbf{n}_d]$ . Theorem 23 below states that each element  $Y_{i_1, \dots, i_d}$  can be represented as a function of a uniform random variable  $\eta_{i_1, \dots, i_d} \perp\!\!\!\perp \mathbf{X}_{\mathbf{n}_d}$ , and of a sequence of canonical forms: one for each sub-array of  $\mathbf{X}_{\mathbf{n}_d}$  that contains  $X_{i_1, \dots, i_d}$ . As was the case for  $d = 2$ , we assume

$$Y_{i_1, \dots, i_d} \perp\!\!\!\perp_{\mathbf{X}_{\mathbf{n}_d}} (\mathbf{Y}_{\mathbf{n}_d} \setminus Y_{i_1, \dots, i_d}), \quad \text{for each } (i_1, \dots, i_d) \in [\mathbf{n}_d]. \quad (61)$$

For convenience, denote by  $\mathbf{i} := (i_1, \dots, i_d)$  an element of the index set  $[\mathbf{n}_d]$ . For each  $k_1 \in [d]$ , let  $\mathbf{X}_{\mathbf{i} \setminus \{k_1\}}^\downarrow$  be the  $(d-1)$ -dimensional sub-array of  $\mathbf{X}_{\mathbf{n}_d}$  obtained by fixing the  $k_1$ th element of  $\mathbf{i}$  and letting all other indices vary over their entire range. For example, for an array with  $d = 3$ ,  $\mathbf{X}_{\mathbf{i} \setminus \{1\}}^\downarrow$  is the matrix extracted from  $\mathbf{X}_{\mathbf{n}_d}$  by fixing  $i_1$ ,  $\mathbf{X}_{i_1, \cdot, \cdot}$ .

Iterating, let  $\mathbf{X}_{\mathbf{i} \setminus \{k_1, \dots, k_p\}}^\downarrow$  be the  $(d-p)$ -dimensional sub-array of  $\mathbf{X}_{\mathbf{n}_d}$  obtained by fixing elements  $i_{k_1}, \dots, i_{k_p}$  of  $\mathbf{i}$  and letting all other indices vary. For each  $p \in [d]$ , denote by  $[d]^p$  the collections of subsets of  $[d]$  with exactly  $p$  elements; let  $[d]^{(p)}$  be the collection of subsets of  $[d]$  with  $p$  or fewer elements. Let the collection of  $(d-p)$ -dimensional sub-arrays containing  $\mathbf{i}$  be denoted as

$$\mathbf{X}_{\mathbf{i}}^{\downarrow(d-p)} = (\mathbf{X}_{\mathbf{i}(s)}^\downarrow)_{s \in [d]^p}.$$

A  $d$ -dimensional version of the augmented canonical form (36) is needed. To that end, let  $\mathbf{j} = (j_1, \dots, j_d) \in [\mathbf{n}_d]$  and define the index sequence  $(\mathbf{i}(s), \mathbf{j}) \in [\mathbf{n}_d]$  as

$$(\mathbf{i}(s), \mathbf{j})_k = \begin{cases} i_k & \text{if } k \in s \\ j_k & \text{if } k \notin s \end{cases}, \quad \text{for } s \in [d]^{(p)} \text{ and } k \in [d]. \quad (62)$$

Then we define the  $d$ -dimensional  $p$ -augmented canonical form,

$$[\mathbf{C}_{\mathbf{i}, \mathbf{X}_{\mathbf{n}_d}}^{(p)}]_{\mathbf{j}} = ([\mathbf{C}_{\mathbf{X}_{\mathbf{n}_d}}]_{\mathbf{j}}, ([\mathbf{C}_{\mathbf{X}_{\mathbf{n}_d}}]_{(\mathbf{i}(s), \mathbf{j})})_{s \in [d]^{(p)}}). \quad (63)$$

Denote by  $\mathcal{C}_{\mathbf{n}_d}^{(d,p)}$  the space of all such augmented canonical forms.

The function returns the collection of elements from  $\mathbf{C}_{\mathbf{X}_{\mathbf{n}_d}}$  that correspond to  $\mathbf{j}$  in each of the  $(d - q)$ -dimensional sub-arrays containing  $\mathbf{i}$ ,  $\mathbf{X}_{\mathbf{i}}^{\downarrow(d-q)}$ , for  $q = 0, \dots, p$ . Alternatively, observe that the function can be constructed by recursing through the  $(d - q)$ -dimensional 1-augmented canonical forms (there are  $\binom{d}{q}$  of them for each  $q$ ), for  $q = 0, \dots, p - 1$ , and returning the first argument of each. This recursive structure captures the action of  $\mathbb{S}_{\mathbf{n}_d}$  on  $\mathbf{X}_{\mathbf{n}_d}$ , and is at the heart of the following result, which gives a functional representation of  $\mathbb{S}_{\mathbf{n}_d}$ -equivariant conditional distributions.

**Theorem 23** *Suppose  $\mathbf{X}_{\mathbf{n}_d}$  and  $\mathbf{Y}_{\mathbf{n}_d}$  are  $\mathcal{X}$ -valued arrays indexed by  $[\mathbf{n}_d]$ , and that  $\mathbf{X}_{\mathbf{n}_d}$  is separately exchangeable. Assume that the elements of  $\mathbf{Y}_{\mathbf{n}_d}$  are conditionally independent given  $\mathbf{X}_{\mathbf{n}_d}$ . Then  $P_{\mathbf{Y}_{\mathbf{n}_d} | \mathbf{X}_{\mathbf{n}_d}}$  is  $\mathbb{S}_{\mathbf{n}_d}$ -equivariant if and only if*

$$(\mathbf{X}_{\mathbf{n}_d}, \mathbf{Y}_{\mathbf{n}_d}) \stackrel{\text{a.s.}}{=} \left( \mathbf{X}_{\mathbf{n}_d}, (f(\eta_{\mathbf{i}}, X_{\mathbf{i}}, \mathbf{C}_{\mathbf{i}, \mathbf{X}_{\mathbf{n}_d}}^{(d)})_{\mathbf{i} \in [\mathbf{n}_d]}) \right), \quad (64)$$

for some measurable function  $f : [0, 1] \times \mathcal{X} \times \mathcal{C}_{\mathbf{n}_d}^{(d,d)} \rightarrow \mathcal{Y}$  and i.i.d. uniform random variables  $(\eta_{\mathbf{i}})_{\mathbf{i} \in [\mathbf{n}_d]} \perp\!\!\!\perp \mathbf{X}_{\mathbf{n}_d}$ .

**Proof** First, assume that  $\mathbf{Y}_{\mathbf{n}_d}$  is conditionally  $\mathbb{S}_{\mathbf{n}_d}$ -equivariant given  $\mathbf{X}_{\mathbf{n}_d}$ . By assumption,  $\mathbf{X}_{\mathbf{n}_d}$  is separately exchangeable, so by Theorem 1,  $\pi_d \cdot (\mathbf{X}_{\mathbf{n}_d}, \mathbf{Y}_{\mathbf{n}_d}) \stackrel{\text{d}}{=} (\mathbf{X}_{\mathbf{n}_d}, \mathbf{Y}_{\mathbf{n}_d})$  for all  $\pi_d \in \mathbb{S}_{\mathbf{n}_d}$ . Fix  $\mathbf{i} \in [\mathbf{n}_d]$ , and let  $\mathbb{S}_{\mathbf{n}_d}^{(\mathbf{i})} \subset \mathbb{S}_{\mathbf{n}_d}$  be the stabilizer of  $\mathbf{i}$ . Observe that each  $\pi_d^{(\mathbf{i})} \in \mathbb{S}_{\mathbf{n}_d}^{(\mathbf{i})}$  fixes  $X_{\mathbf{i}}$  and  $Y_{\mathbf{i}}$ .

In analogy to the fixed -row and -column structure for the  $d = 2$  case, any  $\pi_d^{(\mathbf{i})} \in \mathbb{S}_{\mathbf{n}_d}^{(\mathbf{i})}$  results in sub-arrays of  $\mathbf{X}_{\mathbf{n}_d}$  in which the elements may be rearranged, but the sub-array maintains its position within  $\mathbf{X}_{\mathbf{n}_d}$ . For example, consider  $d = 3$ . Each of the two-dimensional arrays  $\mathbf{X}_{i_1, :, :}$ ,  $\mathbf{X}_{:, i_2, :}$ , and  $\mathbf{X}_{:, :, i_3}$  may have their elements rearranged, but they will remain the two-dimensional sub-arrays that intersect at  $\mathbf{i}$ . Likewise for the one-dimensional sub-arrays  $\mathbf{X}_{i_1, i_2, :}$ ,  $\mathbf{X}_{i_1, :, i_3}$ , and  $\mathbf{X}_{:, i_2, i_3}$ .

Recall that  $\mathbf{X}_{\mathbf{i}}^{\downarrow(\{k_1, \dots, k_p\})}$  is the  $(d - p)$ -dimensional sub-array of  $\mathbf{X}_{\mathbf{n}_d}$  obtained by fixing elements  $i_{k_1}, \dots, i_{k_p}$  of  $\mathbf{i}$  and letting the other indices vary, and  $\mathbf{X}_{\mathbf{i}}^{\downarrow(d-p)}$  is the collection of these sub-arrays. For  $p = 1$ ,  $\mathbf{X}_{\mathbf{i}}^{\downarrow(d-1)} = (\mathbf{X}_{\mathbf{i}(\{k_1\})}^{\downarrow})_{k_1 \in [d]}$  is the collection of  $(d - 1)$ -dimensional sub-arrays containing the element  $X_{\mathbf{i}}$ . Denote by  $\mathbf{X}_{\mathbf{i}}^{\uparrow(d)} := \mathbf{X}_{\mathbf{n}_d} \setminus \mathbf{X}_{\mathbf{i}}^{\downarrow(d-1)}$  the  $d$ -dimensional array that remains after extracting  $\mathbf{X}_{\mathbf{i}}^{\downarrow(d-1)}$  from  $\mathbf{X}_{\mathbf{n}_d}$ . Call  $\mathbf{X}_{\mathbf{i}}^{\uparrow(d)}$  a  $d$ -dimensional *remainder*

array. For example, with  $d = 3$ , the remainder array  $\mathbf{X}_i^{\uparrow(d)}$  consists of  $\mathbf{X}_{\mathbf{n}_d}$  with the three two-dimensional arrays (matrices) that contain  $X_i$  removed.

Continue recursively and construct a remainder array  $\mathbf{X}_{i(s)}^{\uparrow(d-p)}$ ,  $s \in [d]^p$ , from each subarray in the collection  $\mathbf{X}_i^{\downarrow(d-p)} = (\mathbf{X}_{i(s)}^{\downarrow})_{s \in [d]^p}$ . The recursive structure of the collection of remainder arrays  $((\mathbf{X}_{i(s)}^{\uparrow(d-p)})_{s \in [d]^p})_{0 \leq p \leq d-1}$  is captured by the array  $\mathbf{Z}^{(i)}$ , with entries

$$[\mathbf{Z}^{(i)}]_{\mathbf{j}} = (X_{\mathbf{j}}, (X_{(\mathbf{i}(s), \mathbf{j})})_{s \in [d]^{(d)}}), \quad (65)$$

where  $(\mathbf{i}(s), \mathbf{j})$  is as in (62). Define the action of a collection of permutations  $\pi'_d \in \mathbb{S}_{\mathbf{n}_d - \mathbf{1}_d}$  on  $\mathbf{Z}^{(i)}$  to be such that, with  $\mathbf{j}_\pi = (\pi_1(j_1), \dots, \pi_d(j_d))$ ,

$$[\pi'_d \cdot \mathbf{Z}^{(i)}]_{\mathbf{j}} = (X_{\mathbf{j}_\pi}, (X_{(\mathbf{i}(s), \mathbf{j}_\pi)})_{s \in [d]^{(d)}}). \quad (66)$$

$\mathbf{Z}^{(i)}$  inherits the exchangeability of  $\mathbf{X}_{\mathbf{n}_d}$ , so that marginally for  $Y_i$ ,

$$(\pi_d \cdot \mathbf{Z}^{(i)}, (X_i, Y_i)) \stackrel{d}{=} (\mathbf{Z}^{(i)}, (X_i, Y_i)) \quad \text{for all } \pi_d \in \mathbb{S}_{\mathbf{n}_d - \mathbf{1}_d}.$$

which implies  $(X_i, Y_i) \perp\!\!\!\perp_{\mathbf{C}_{\mathbf{Z}^{(i)}}} \mathbf{Z}^{(i)}$ . Conditioning on  $X_i$  and  $\mathbf{C}_{\mathbf{Z}^{(i)}}$  is the same as conditioning on  $X_i$  and the  $d$ -dimensional  $d$ -augmented canonical form  $\mathbf{C}_{i, \mathbf{X}_{\mathbf{n}_d}}^{(d)}$  defined in (63), implying that

$$Y_i \perp\!\!\!\perp_{(X_i, \mathbf{C}_{i, \mathbf{X}_{\mathbf{n}_d}}^{(d)})} \mathbf{X}_{\mathbf{n}_d}.$$

By Theorem 5, there is a measurable function  $f_i : [0, 1] \times \mathcal{X} \times \mathcal{C}_{\mathbf{n}_d}^{(d)} \rightarrow \mathcal{Y}$  such that

$$Y_i = f_i(\eta_i, X_i, \mathbf{C}_{i, \mathbf{X}_{\mathbf{n}_d}}^{(d)}),$$

for a uniform random variable  $\eta_i \perp\!\!\!\perp \mathbf{X}_{\mathbf{n}_d}$ . This is true for all  $\mathbf{i} \in [\mathbf{n}_d]$ ; by equivariance the same  $f_i = f$  must work for every  $\mathbf{i}$ . Furthermore, by assumption the elements of  $\mathbf{Y}_{\mathbf{n}_d}$  are mutually conditionally independent given  $\mathbf{X}_{\mathbf{n}_d}$ , and therefore by the chain rule for conditional independence (Kallenberg, 2002, Prop. 6.8), the joint identity (64) holds.

The converse is straightforward to verify. ■

## References

M. Abadi, A. Agarwal, P. Barham, E. Brevdo, Z. Chen, C. Citro, G. S. Corrado, A. Davis, J. Dean, M. Devin, S. Ghemawat, I. Goodfellow, A. Harp, G. Irving, M. Isard, Y. Jia, R. Jozefowicz, L. Kaiser, M. Kudlur, J. Levenberg, D. Mané, R. Monga, S. Moore, D. Murray, C. Olah, M. Schuster, J. Shlens, B. Steiner, I. Sutskever, K. Talwar, P. Tucker, V. Vanhoucke, V. Vasudevan, F. Viégas, O. Vinyals, P. Warden, M. Wattemberg, M. Wicke, Y. Yu, and X. Zheng. TensorFlow: Large-scale machine learning on heterogeneous systems, 2015. URL <https://www.tensorflow.org/>. Software available from tensorflow.org.

- D. J. Aldous. Exchangeability and related topics. In P. L. Hennequin, editor, *École d'Été de Probabilités de Saint-Flour XIII - 1983*, number 1117 in Lecture Notes in Mathematics, pages 1–198. Springer, 1985.
- S. Andersson. Distributions of maximal invariants using quotient measures. *The Annals of Statistics*, 10(3):955–961, Sep 1982.
- S. Arora, R. Ge, B. Neyshabur, and Y. Zhang. Stronger generalization bounds for deep nets via a compression approach. In *International Conference on Machine Learning (ICML)*, pages 254–263, 2018.
- T. Austin. Exchangeable random measures. *Annales de l'Institut Henri Poincaré, Probabilités et Statistiques*, 51(3):842–861, 08 2015.
- S. Bacallado, S. Favaro, and L. Trippa. Bayesian nonparametric analysis of reversible Markov chains. *The Annals of Statistics*, 41(2):870–896, 2013.
- R. R. Bahadur. Sufficiency and statistical decision functions. *The Annals of Mathematical Statistics*, 25(3):423–462, 1954.
- I. Bello, B. Zoph, A. Vaswani, J. Shlens, and Q. V. Le. Attention augmented convolutional networks. In *The IEEE International Conference on Computer Vision (ICCV)*, 2019.
- A. Bietti and J. Mairal. Group invariance, stability to deformations, and complexity of deep convolutional representations. *Journal of Machine Learning Research*, 20(25):1–49, 2019.
- D. Blackwell and J. B. MacQueen. Ferguson distributions via polya urn schemes. *The Annals of Statistics*, 1(2):353–355, Mar 1973.
- B. Bloem-Reddy and Y. W. Teh. Neural network models of exchangeable sequences. NeurIPS Workshop on Bayesian Deep Learning, 2018.
- C. Borgs and J. Chayes. Graphons: A nonparametric method to model, estimate, and design algorithms for massive networks. *Proceedings of the 2017 ACM Conference on Economics and Computation - EC '17*, 2017.
- J. Bruna and S. Mallat. Invariant scattering convolution networks. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 35(8):1872–1876, 2013.
- J. Bruna and S. Mallat. Multiscale sparse microcanonical models. *Mathematical Statistics and Learning*, 1(3/4):257–315, 2018.
- J. Bruna, W. Zaremba, A. Szlam, and Y. LeCun. Spectral networks and locally connected networks on graphs. In *International Conference on Learning Representations (ICLR)*, 2014.
- D. Cai, T. Campbell, and T. Broderick. Edge-exchangeable graphs and sparsity. In *Advances in Neural Information Processing Systems (Neurips)*, pages 4249–4257. 2016.

- F. Caron and E. B. Fox. Sparse graphs using exchangeable random measures. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 79(5):1–44, 2017.
- E. Çinlar. *Probability and Stochastics*. Springer New York, 2011.
- J. Chan, V. Perrone, J. Spence, P. Jenkins, S. Mathieson, and Y. Song. A likelihood-free inference framework for population genetic data using exchangeable neural networks. In *Advances in Neural Information Processing Systems (Neurips)*, pages 8594–8605. 2018.
- S. Chen, E. Dobriban, and J. H. Lee. Invariance reduces variance: Understanding data augmentation in deep learning and beyond. *arXiv e-prints*, abs/1907.10905, 2019.
- T. S. Cohen and M. Welling. Group equivariant convolutional networks. In *International Conference on Machine Learning (ICML)*, pages 2990–2999, 2016.
- T. S. Cohen and M. Welling. Steerable CNNs. In *International Conference on Learning Representations (ICLR)*, 2017.
- T. S. Cohen, M. Geiger, J. Köhler, and M. Welling. Spherical CNNs. In *International Conference on Learning Representations (ICLR)*, 2018.
- T. S. Cohen, M. Geiger, and M. Weiler. A general theory of equivariant cnns on homogeneous spaces. In *Advances in Neural Information Processing Systems (Neurips)*, pages 9145–9156. 2019.
- D. R. Cox and D. V. Hinkley. *Theoretical Statistics*. Chapman & Hall, London, 1974.
- H. Crane and W. Dempsey. Edge exchangeable models for interaction networks. *Journal of the American Statistical Association*, 113(523):1311–1326, 2018.
- G. Cybenko. Approximation by superpositions of a sigmoidal function. *Mathematics of Control, Signals, and Systems*, 2(4):303–314, Dec 1989.
- A. Dawid. Invariance and independence in multivariate distribution theory. *Journal of Multivariate Analysis*, 17(3):304–315, 1985.
- B. de Finetti. Fuzione caratteristica di un fenomeno aleatorio. *Mem. R. Acc. Lincei*, 6(4):86–133, 1930.
- P. Diaconis. Finite forms of de finetti’s theorem on exchangeability. *Synthese*, 36(2):271–281, 1977.
- P. Diaconis. Sufficiency as statistical symmetry. In F. Browder, editor, *Proceedings of the AMS Centennial Symposium*, pages 15–26. American Mathematical Society, 1988.
- P. Diaconis and D. Freedman. Finite exchangeable sequences. *The Annals of Probability*, 8(4):745–764, 1980a.
- P. Diaconis and D. Freedman. De Finetti’s Theorem for Markov Chains. *The Annals of Probability*, 8(1):115–130, 02 1980b.

- P. Diaconis and D. Freedman. Partial exchangeability and sufficiency. In J. K. Ghosh and J. Roy, editors, *Proc. Indian Stat. Inst. Golden Jubilee Int'l Conf. Stat.: Applications and New Directions*, pages 205–236. Indian Statistical Institute, 1984.
- P. Diaconis and D. Freedman. A dozen de Finetti-style results in search of a theory. *Annales de l'Institut Henri Poincaré, Probabilités et Statistiques*, 23:397–423, 1987.
- D. K. Duvenaud, D. Maclaurin, J. Iparraguirre, R. Bombarell, T. Hirzel, A. Aspuru-Guzik, and R. P. Adams. Convolutional networks on graphs for learning molecular fingerprints. In *Advances in Neural Information Processing Systems (Neurips)*, pages 2224–2232. 2015.
- M. L. Eaton. *Group invariance in applications in statistics*, volume 1 of *Regional Conference Series in Probability and Statistics*. Institute of Mathematical Statistics and American Statistical Association, Haywood, CA and Alexandria, VA, 1989.
- H. Edwards and A. Storkey. Towards a neural statistician. In *International Conference on Learning Representations (ICLR)*, 2017.
- R. A. Fisher. On the mathematical foundations of theoretical statistics. *Philosophical Transactions of the Royal Society A: Mathematical, Physical and Engineering Sciences*, 222(594-604):309—368, 1922.
- S. Fortini, L. Ladelli, and E. Regazzini. Exchangeability, predictive distributions and parametric models. *Sankhyā: The Indian Journal of Statistics*, 62(1):86–109, 2000.
- D. A. Freedman. Invariants under mixing which generalize de Finetti's theorem. *The Annals of Mathematical Statistics*, 33(3):916–923, 1962.
- D. A. Freedman. Invariants under mixing which generalize de Finetti's theorem: Continuous time parameter. *The Annals of Mathematical Statistics*, 34(4):1194–1216, 1963.
- F. Gao, G. Wolf, and M. Hirn. Geometric scattering for graph data analysis. In *International Conference on Machine Learning (ICML)*, pages 2122–2131, 2019.
- M. Garnelo, D. Rosenbaum, C. Maddison, T. Ramalho, D. Saxton, M. Shanahan, Y. W. Teh, D. Rezende, and S. M. A. Eslami. Conditional neural processes. In *International Conference on Machine Learning (ICML)*, pages 1704–1713, 2018.
- R. Gens and P. M. Domingos. Deep symmetry networks. In *Advances in Neural Information Processing Systems (Neurips)*, pages 2537–2545. 2014.
- J. Gilmer, S. S. Schoenholz, P. F. Riley, O. Vinyals, and G. E. Dahl. Neural message passing for quantum chemistry. In *International Conference on Machine Learning (ICML)*, pages 1263–1272, 2017.
- N. C. Giri. *Group Invariance in Statistical Inference*. World Scientific, 1996.
- I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio. Generative adversarial nets. In *Advances in Neural Information Processing Systems (Neurips)*, pages 2672–2680. 2014.

- W. Grathwohl, K.-C. Wang, J.-H. Jacobsen, D. Duvenaud, M. Norouzi, and K. Swersky. Your classifier is secretly an energy based model and you should treat it like one. In *International Conference on Learning Representations (ICLR)*, 2020.
- M. U. Gutmann and J. Corander. Bayesian optimization for likelihood-free inference of simulator-based statistical models. *Journal of Machine Learning Research*, 17(125):1–47, 2016.
- W. J. Hall, R. A. Wijsman, and J. K. Ghosh. The relationship between sufficiency and invariance with applications in sequential analysis. *The Annals of Mathematical Statistics*, 36(2):575–614, 1965. ISSN 00034851.
- P. R. Halmos and L. J. Savage. Application of the Radon-Nikodym theorem to the theory of sufficient statistics. *The Annals of Mathematical Statistics*, 20(2):225–241, 1949.
- T. Hamaguchi, H. Oiwa, M. Shimbo, and Y. Matsumoto. Knowledge transfer for out-of-knowledge-base entities: A graph neural network approach. In *International Joint Conference on Artificial Intelligence (IJCAI)*, pages 1802–1808, 2017.
- J. Hartford, D. Graham, K. Leyton-Brown, and S. Ravanbakhsh. Deep models of interactions across sets. In *International Conference on Machine Learning (ICML)*, pages 1914–1923, 2018.
- K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 770–778, 2016.
- R. Herzig, M. Raboh, G. Chechik, J. Berant, and A. Globerson. Mapping images to scene graphs with permutation-invariant structured prediction. In *Advances in Neural Information Processing Systems (Neurips)*, pages 7211–7221. 2018.
- E. Hewitt and L. J. Savage. Symmetric measures on Cartesian products. *Transactions of the American Mathematical Society*, 80(2):470–501, 1955.
- G. Hinton, S. Sabour, and N. Frosst. Matrix capsules with em routing. In *International Conference on Learning Representations (ICLR)*, 2018.
- K. Hornik. Approximation capabilities of multilayer feedforward networks. *Neural Networks*, 4(2):251–257, 1991.
- K. Hornik, M. Stinchcombe, and H. White. Multilayer feedforward networks are universal approximators. *Neural Networks*, 2(5):359–366, 1989.
- E. T. Jaynes. Information theory and statistical mechanics. *Phys. Rev.*, 106:620–630, May 1957.
- O. Kallenberg. *Foundations of Modern Probability*. Springer-Verlag New York, 2nd edition, 2002.
- O. Kallenberg. *Probabilistic Symmetries and Invariance Principles*. Springer, 2005.



- O. Kallenberg. *Random Measures, Theory and Applications*. Springer International Publishing, 2017.
- S. Kearnes, K. McCloskey, M. Berndl, V. Pande, and P. Riley. Molecular graph convolutions: moving beyond fingerprints. *Journal of Computer-Aided Molecular Design*, 30(8):595–608, Aug 2016.
- N. Keriven and G. Peyré. Universal invariant and equivariant graph neural networks. In *Advances in Neural Information Processing Systems (Neurips)*, pages 7092–7101. 2019.
- H. Kim, A. Mnih, J. Schwarz, M. Garnelo, A. Eslami, D. Rosenbaum, O. Vinyals, and Y. W. Teh. Attentive neural processes. In *International Conference on Learning Representations (ICLR)*, 2019.
- D. P. Kingma and M. Welling. Auto-Encoding Variational Bayes. In *Proceedings of the 2nd International Conference on Learning Representations (ICLR)*, 2014.
- R. Kondor and S. Trivedi. On the generalization of equivariance and convolution in neural networks to the action of compact groups. In *International Conference on Machine Learning (ICML)*, pages 2747–2755, 2018.
- I. Korshunova, J. Degraeve, F. Huszar, Y. Gal, A. Gretton, and J. Dambre. BRUNO: A deep recurrent model for exchangeable data. In *Advances in Neural Information Processing Systems (Neurips)*, pages 7190–7198. 2018.
- S. L. Lauritzen. Sufficiency, prediction and extreme models. *Scandinavian Journal of Statistics*, 1(3):128–134, 1974a.
- S. L. Lauritzen. On the interrelationships among sufficiency, total sufficiency and some related concepts. Technical Report 8, Institute of Mathematical Statistics, University of Copenhagen, July 1974b.
- S. L. Lauritzen. Extreme point models in statistics (with discussion and reply). *Scandinavian Journal of Statistics*, 11(2):65–91, 1984.
- S. L. Lauritzen. *Extremal Families and Systems of Sufficient Statistics*. Lecture Notes in Statistics. Springer, 1988.
- S. L. Lauritzen. *Graphical Models*. Oxford University Press, 1996.
- Y. LeCun, B. Boser, J. S. Denker, D. Henderson, R. E. Howard, W. Hubbard, and L. D. Jackel. Backpropagation applied to handwritten zip code recognition. *Neural Computation*, 1(4):541–551, 1989.
- Y. LeCun, S. Chopra, R. Hadsell, M. Ranzato, and F.-J. Huang. *Predicting Structured Data*, chapter A Tutorial on Energy-Based Learning. 2006.
- J. Lee, Y. Lee, J. Kim, A. Kosiorek, S. Choi, and Y. W. Teh. Set transformer: A framework for attention-based permutation-invariant neural networks. In *International Conference on Machine Learning (ICML)*, pages 3744–3753, 2019.

- E. L. Lehmann and J. P. Romano. *Testing Statistical Hypotheses*. Springer Texts in Statistics. Springer-Verlag New York, 2005.
- J. E. Lenssen, M. Fey, and P. Libuschewski. Group equivariant capsule networks. In *Advances in Neural Information Processing Systems (Neurips)*, pages 8844–8853. 2018.
- L. Lovász. *Large Networks and Graph Limits*. American Mathematical Society, 2012.
- C. Lyle, M. van der Wilk, M. Kwiatkowska, Y. Gal, and B. Bloem-Reddy. On the benefits of invariance in neural networks. 2020.
- H. Mahmoud. *Pólya Urn Models*. Chapman & Hall/CRC Texts in Statistical Science. Chapman & Hall/CRC, 2008.
- H. Maron, E. Fetaya, N. Segol, and Y. Lipman. On the universality of invariant networks. In *International Conference on Machine Learning (ICML)*, pages 4363–4371, 2019.
- B. D. McKay and A. Piperno. Practical graph isomorphism, II. *Journal of Symbolic Computation*, 60:94 – 112, 2014.
- M. L. Minsky and S. A. Papert. *Perceptrons: Expanded Edition*. Cambridge, MA, USA, 1988.
- R. L. Murphy, B. Srinivasan, V. Rao, and B. Ribeiro. Janosy pooling: Learning deep permutation-invariant functions for variable-size inputs. In *International Conference on Learning Representations (ICLR)*, 2019.
- M. Niepert and G. Van den Broeck. Tractability through exchangeability: A new perspective on efficient probabilistic inference. In *AAAI Conference on Artificial Intelligence*, pages 2467–2475, 2014.
- M. Niepert, M. Ahmed, and K. Kutzkov. Learning convolutional neural networks for graphs. In *International Conference on Machine Learning (ICML)*, pages 2014–2023, 2016.
- P. Orbanz and D. M. Roy. Bayesian models of graphs, arrays and other exchangeable random structures. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 37(2):437–461, Feb 2015.
- N. Parmar, P. Ramachandran, A. Vaswani, I. Bello, A. Levskaya, and J. Shlens. Stand-alone self-attention in vision models. In *Advances in Neural Information Processing Systems (Neurips)*, pages 68–80. 2019.
- A. Paszke, S. Gross, F. Massa, A. Lerer, J. Bradbury, G. Chanan, T. Killeen, Z. Lin, N. Gimelshein, L. Antiga, A. Desmaison, A. Kopf, E. Yang, Z. DeVito, M. Raison, A. Tejani, S. Chilamkurthy, B. Steiner, L. Fang, J. Bai, and S. Chintala. PyTorch: An imperative style, high-performance deep learning library. In *Advances in Neural Information Processing Systems (Neurips)*, pages 8026–8037. 2019.
- S. Ravanbakhsh. Universal equivariant multilayer perceptrons. *arXiv e-prints*, abs/2002.02912, 2020.

- S. Ravanbakhsh, J. Schneider, and B. Póczos. Equivariance through parameter-sharing. In *International Conference on Machine Learning (ICML)*, pages 2892–2901, 2017.
- D. J. Rezende, S. Mohamed, and D. Wierstra. Stochastic backpropagation and approximate inference in deep generative models. In *International Conference on Machine Learning (ICML)*, number 2, pages 1278–1286, 2014.
- D. W. Romero, E. J. Bekkers, J. M. Tomczak, and M. Hoogendoorn. Attentive group equivariant convolutional networks. *arXiv e-prints*, abs/2002.03830, 2020.
- J. Rotman. *An Introduction to the Theory of Groups*, volume 148 of *Graduate Texts in Mathematics*. Springer-Verlag New York, 4 edition, 1995.
- S. Sabour, N. Frosst, and G. E. Hinton. Dynamic routing between capsules. In *Advances in Neural Information Processing Systems (Neurips)*, pages 3856–3866. 2017.
- F. Scarselli, M. Gori, A. C. Tsoi, M. Hagenbuchner, and G. Monfardini. The graph neural network model. *IEEE Transactions on Neural Networks*, 20(1):61–80, 2009.
- M. J. Schervish. *Theory of Statistics*. Springer-Verlag New York, 1995.
- W. Schindler. *Measures with Symmetry Properties*, volume 1808 of *Lecture Notes in Mathematics*. Springer-Verlag Berlin Heidelberg, Berlin, 2003.
- N. Segol and Y. Lipman. On universal equivariant set networks. In *International Conference on Learning Representations (ICLR)*, 2020.
- J. Shawe-Taylor. Building symmetries into feedforward networks. In *1989 First IEE International Conference on Artificial Neural Networks, (Conf. Publ. No. 313)*, pages 158–162, Oct 1989.
- J. Shawe-Taylor. Threshold network learning in the presence of equivalences. In *Advances in Neural Information Processing Systems (Neurips)*, pages 879–886. Morgan-Kaufmann, 1991.
- J. Shawe-Taylor. Symmetries and discriminability in feedforward network architectures. *IEEE Transactions on Neural Networks*, 4(5):816–826, 1993.
- J. Shawe-Taylor. Sample sizes for threshold networks with equivalences. *Information and Computation*, 118(1):65 – 72, 1995.
- M. Skibinsky. Adequate subfields and sufficiency. *The Annals of Mathematical Statistics*, 38(1):155–161, 1967. ISSN 00034851.
- P. Smolensky. *Information Processing in Dynamical Systems: Foundations of Harmony Theory*, pages 194–281. 1987.
- E. Snelson and Z. Ghahramani. Sparse gaussian processes using pseudo-inputs. In *Advances in Neural Information Processing Systems (Neurips)*, pages 1257–1264. 2006.

- T. P. Speed. A factorisation theorem for adequate statistics. *Australian Journal of Statistics*, 20(3):240–249, 1978.
- N. Srivastava, G. Hinton, A. Krizhevsky, I. Sutskever, and R. Salakhutdinov. Dropout: A simple way to prevent neural networks from overfitting. *Journal of Machine Learning Research*, 15:1929–1958, 2014.
- S. Sukhbaatar, A. Szlam, J. Weston, and R. Fergus. End-to-end memory networks. In *Advances in Neural Information Processing Systems (Neurips)*, pages 2440–2448. 2015.
- The GAP Group. GAP – Groups, Algorithms, and Programming, Version 4.10.0, 2018. URL <https://www.gap-system.org>.
- Theano Development Team. Theano: A Python framework for fast computation of mathematical expressions. *arXiv e-prints*, abs/1605.02688, May 2016.
- J.-W. van de Meent, B. Paige, H. Yang, and F. Wood. An introduction to probabilistic programming. 09 2018. URL <https://arxiv.org/pdf/1809.10756>.
- A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. u. Kaiser, and I. Polosukhin. Attention is all you need. In *Advances in Neural Information Processing Systems (Neurips)*, pages 5998–6008. 2017.
- S. Wager, W. Fithian, S. Wang, and P. S. Liang. Altitude training: Strong bounds for single-layer dropout. In *Advances in Neural Information Processing Systems (Neurips)*, pages 100–108. 2014.
- E. Wagstaff, F. Fuchs, M. Engelcke, I. Posner, and M. A. Osborne. On the limitations of representing functions on sets. In *International Conference on Machine Learning (ICML)*, pages 6487–6494, 2019.
- M. Welling, M. Rosen-zvi, and G. E. Hinton. Exponential family harmoniums with an application to information retrieval. In *Advances in Neural Information Processing Systems (Neurips)*, pages 1481–1488. 2005.
- R. A. Wijsman. *Invariant measures on groups and their use in statistics*, volume 14 of *Lecture Notes–Monograph Series*. Institute of Mathematical Statistics, Hayward, CA, 1990.
- S. A. Williamson. Nonparametric network models for link prediction. *Journal of Machine Learning Research*, 17(202):1–21, 2016.
- S. Wiqvist, P.-A. Mattei, U. Picchini, and J. Frellsen. Partially exchangeable networks and architectures for learning summary statistics in approximate Bayesian computation. In *International Conference on Machine Learning (ICML)*, pages 6798–6807, 2019.
- Wolfram Research, Inc. Mathematica, Version 11.3, 2018. Champaign, IL.
- J. Wood and J. Shawe-Taylor. Representation theory and invariant neural networks. *Discrete Applied Mathematics*, 69(1):33–60, 1996.

- K. Xu, W. Hu, J. Leskovec, and S. Jegelka. How powerful are graph neural networks? In *International Conference on Learning Representations (ICLR)*, 2019.
- D. Yarotsky. Universal approximations of invariant maps by neural networks. *arXiv e-prints*, abs/1804.10306, 2018.
- F. Yu and V. Koltun. Multi-scale context aggregation by dilated convolutions. In *International Conference on Learning Representations (ICLR)*, 2016.
- S. L. Zabell. *The Rule of Succession*, pages 38–73. Cambridge Studies in Probability, Induction and Decision Theory. Cambridge University Press, 2005.
- M. Zaheer, S. Kottur, S. Ravanbakhsh, B. Poczos, R. R. Salakhutdinov, and A. J. Smola. Deep sets. In *Advances in Neural Information Processing Systems (Neurips)*, pages 3391–3401. 2017.
- W. Zhou, V. Veitch, M. Austern, R. P. Adams, and P. Orbanz. Non-vacuous generalization bounds at the ImageNet scale: a PAC-Bayesian compression approach. In *International Conference on Learning Representations (ICLR)*, 2019.