

# Fourier Theoretic Probabilistic Inference over Permutations

**Jonathan Huang**

*Robotics Institute  
Carnegie Mellon University  
Pittsburgh, PA 15213*

JCH1@CS.CMU.EDU

**Carlos Guestrin**

*Machine Learning Department  
Carnegie Mellon University  
Pittsburgh, PA 15213*

GUESTRIN@CS.CMU.EDU

**Leonidas Guibas**

*Department of Computer Science  
Stanford University  
Stanford, CA 94305*

GUIBAS@CS.STANFORD.EDU

**Editor:** Marina Meila

## Abstract

Permutations are ubiquitous in many real-world problems, such as voting, ranking, and data association. Representing uncertainty over permutations is challenging, since there are  $n!$  possibilities, and typical compact and factorized probability distribution representations, such as graphical models, cannot capture the mutual exclusivity constraints associated with permutations. In this paper, we use the “low-frequency” terms of a Fourier decomposition to represent distributions over permutations compactly. We present *Kronecker conditioning*, a novel approach for maintaining and updating these distributions directly in the Fourier domain, allowing for polynomial time bandlimited approximations. Low order Fourier-based approximations, however, may lead to functions that do not correspond to valid distributions. To address this problem, we present a quadratic program defined directly in the Fourier domain for projecting the approximation onto a relaxation of the polytope of legal marginal distributions. We demonstrate the effectiveness of our approach on a real camera-based multi-person tracking scenario.

**Keywords:** identity management, permutations, approximate inference, group theoretical methods, sensor networks

## 1. Introduction

Probability distributions over permutations arise in a diverse variety of real world problems. While they were perhaps first studied in the context of gambling and card games, they have now been found to be applicable to many important problems such as multi-object tracking, information retrieval, webpage ranking, preference elicitation, and voting. Probabilistic reasoning problems over permutations, however, are not amenable to the typical representations afforded by machine learning such as Bayesian networks and Markov random fields. This paper explores an alternative representation and inference algorithms based on Fourier analysis for dealing with permutations.

As an example, consider the problem of tracking  $n$  people based on a set of noisy measurements of identity and position. A typical tracking system might attempt to manage a set of  $n$  tracks along

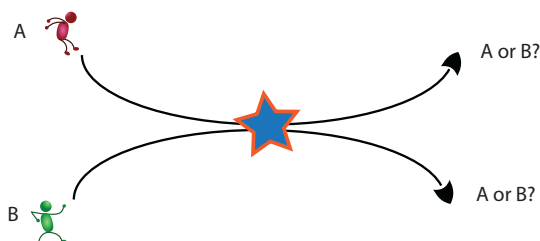


Figure 1: When two persons pass near each other, their identities can get confused.

with an identity corresponding to each track, in spite of ambiguities arising from imperfect identity measurements. When the people are well separated, the problem is easily decomposed and measurements about each individual can be clearly associated with a particular track. When people pass near each other, however, confusion can arise as their signal signatures may mix; see Figure 1. After the individuals separate again, their positions may be clearly distinguishable, but their identities can still be confused, resulting in identity uncertainty which must be propagated forward in time with each person, until additional observations allow for disambiguation. This task of maintaining a belief state for the correct association between object tracks and object identities while accounting for local mixing events and sensor observations, was introduced in Shin et al. (2003) and is called the *identity management problem*.

The identity management problem poses a challenge for probabilistic inference because it needs to address the fundamental combinatorial challenge that there is a factorial number of associations to maintain between tracks and identities. Distributions over the space of all permutations require storing at least  $n! - 1$  numbers, an infeasible task for all but very small  $n$ . Moreover, typical compact representations, such as graphical models, cannot efficiently capture the mutual exclusivity constraints associated with permutations.

While there have been many approaches for coping with the factorial complexity of maintaining a distribution over permutations, most attack the problem using one of two ideas—storing and updating a small subset of likely permutations, or, as in our case, restricting consideration to a tractable subspace of possible distributions. Willsky (1978) was the first to formulate the probabilistic filtering/smoothing problem for group-valued random variables. He proposed an efficient FFT based approach of transforming between primal and Fourier domains so as to avoid costly convolutions, and provided efficient algorithms for dihedral and metacyclic groups. Kueh et al. (1999) show that probability distributions on the group of permutations are well approximated by a small subset of Fourier coefficients of the actual distribution, allowing for a principled tradeoff between accuracy and complexity. The approach taken in Shin et al. (2005), Schumitsch et al. (2005), and Schumitsch et al. (2006) can be seen as an algorithm for maintaining a particular fixed subset of Fourier coefficients of the log density. Most recently, Kondor et al. (2007) allow for a general set of Fourier coefficients, but assume a restrictive form of the observation model in order to exploit an efficient FFT factorization.

In the following, we outline our main contributions and provide a roadmap of the sections ahead.<sup>1</sup>

1. A much shorter version this work appeared in *NIPS 2007* (Huang et al., 2007). We provide a more complete discussion of our Fourier based methods in this extended paper.

- In Sections 4 and 5, we provide a gentle introduction to the theory of group representations and noncommutative Fourier analysis. While none of the results of these sections are novel, and have indeed been studied by mathematicians for decades (Diaconis, 1989; Terras, 1999; Willsky, 1978; Chen, 1989), noncommutative Fourier analysis is still fairly new to the machine learning community, which has just begun to discover some of its exciting applications (Huang et al., 2007, 2009; Kondor et al., 2007; Kondor and Borgwardt, 2008). Our tutorial sections are targeted specifically at the machine learning community and describe its connections to probabilistic inference problems that involve permutations.
- In Section 6, we discuss performing probabilistic inference operations in the Fourier domain. In particular, we present Fourier theoretic algorithms for two ubiquitous operations which appear in filtering applications and beyond: prediction/rollup and conditioning with Bayes rule. Our main contribution in this section is a novel and conceptually simple algorithm, called *Kronecker Conditioning*, which performs all conditioning operations completely in the Fourier domain, allowing for a principled tradeoff between computational complexity and approximation accuracy. Our approach generalizes upon previous work in two ways—first, in the sense that it can address any transition model or likelihood function that can be represented in the Fourier domain, and second, in the sense that many of our results hold for arbitrary finite groups.
- In Section 7, we analyze the errors which can be introduced by bandlimiting a probability distribution and show how they propagate with respect to inference operations. We argue that approximate conditioning based on bandlimited distributions can sometimes yield Fourier coefficients which do not correspond to any valid distribution, even returning negative “probabilities” on occasion. We address possible negative and inconsistent probabilities by presenting a method for projecting the result back into the polytope of coefficients which correspond to nonnegative and consistent marginal probabilities using a simple quadratic program.
- In Section 8, we present a collection of general techniques for efficiently computing the Fourier coefficients of probabilistic models that might be useful in practical inference problems, and give a variety of examples of such computations for probabilistic models that might arise in identity management or ranking scenarios.
- Finally in Section 10, we empirically evaluate the accuracy of approximate inference on simulated data drawn from our model and further demonstrate the effectiveness of our approach on a real camera-based multi-person tracking scenario.

## 2. Filtering Over Permutations

As a prelude to the general problem statement, we begin with a simple identity management problem on three tracks (illustrated in Figure 2) which we will use as a running example. In this problem, we observe a stream of localization data from three people walking inside a room. Except for a camera positioned at the entrance, however, there is no way to distinguish between identities once they are inside. In this example, an internal tracker declares that two tracks have ‘mixed’ whenever they get too close to each other and announces the identity of any track that enters or exits the room.

In our particular example, three people, Alice, Bob and Cathy, enter a room separately, walk around, and we observe Bob as he exits. The events for our particular example in the figure are

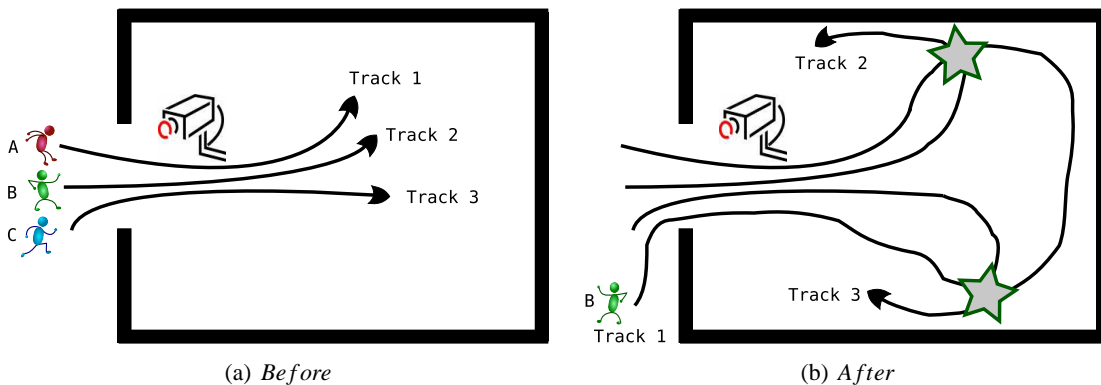


Figure 2: Identity Management example. Three people, Alice, Bob and Charlie enter a room and we receive a position measurement for each person at each time step. With no way to observe identities inside the room, however, we are confused whenever two tracks get too close. In this example, track 1 crosses with track 2, then with track 3, then leaves the room, at which point it is observed that the identity at Track 1 is in fact Bob.

recorded in Table 1. Since Tracks 2 and 3 never mix, we know that Cathy cannot be in Track 2 in the end, and furthermore, since we observe Bob to be in Track 1 when he exits, we can deduce that Cathy must have been in Track 3, and therefore Alice must have been in Track 2. Our simple example illustrates the combinatorial nature of the problem—in particular, reasoning about the mixing events allows us to exactly decide where Alice and Cathy were even though we only made an observation about Bob at the end.

<i>Event #</i>	<i>Event Type</i>
<i>1</i>	<i>Tracks 1 and 2 mixed</i>
<i>2</i>	<i>Tracks 1 and 3 mixed</i>
<i>3</i>	<i>Observed Identity Bob at Track 1</i>

Table 1: Table of Mixing and Observation events logged by the tracker.

In identity management, a permutation  $\sigma$  represents a joint assignment of identities to internal tracks, with  $\sigma(i)$  being the track belonging to the  $i$ th identity. When people walk too closely together, their identities can be confused, leading to uncertainty over  $\sigma$ . To model this uncertainty, we use a *Hidden Markov Model (HMM)* on permutations, which is a joint distribution over latent permutations  $\sigma^{(1)}, \dots, \sigma^{(T)}$ , and observed variables  $z^{(1)}, \dots, z^{(T)}$  which factors as:

$$P(\sigma^{(1)}, \dots, \sigma^{(T)}, z^{(1)}, \dots, z^{(T)}) = P(\sigma^{(1)}) P(z^{(1)} | \sigma^{(1)}) \prod_{t=2}^T P(z^{(t)} | \sigma^{(t)}) \cdot P(\sigma^{(t)} | \sigma^{(t-1)}).$$

The conditional probability distribution  $P(\sigma^{(t)} | \sigma^{(t-1)})$  is called the *transition model*, and might reflect, for example, that the identities belonging to two tracks were swapped with some probability by a mixing event. The distribution  $P(z^{(t)} | \sigma^{(t)})$  is called the *observation model*, which might, for example, capture a distribution over the color of clothing for each individual.

We focus on *filtering*, in which one queries the HMM for the posterior at some time step, conditioned on all past observations. Given the distribution  $P(\sigma^{(t)}|z^{(1)}, \dots, z^{(t)})$ , we recursively compute  $P(\sigma^{(t+1)}|z^{(1)}, \dots, z^{(t+1)})$  in two steps: a *prediction/rollup* step and a *conditioning* step. Taken together, these two steps form the well known *Forward Algorithm* (Rabiner, 1989). The prediction/rollup step multiplies the distribution by the transition model and marginalizes out the previous time step:

$$P(\sigma^{(t+1)}|z^{(1)}, \dots, z^{(t)}) = \sum_{\sigma^{(t)}} P(\sigma^{(t+1)}|\sigma^{(t)})P(\sigma^{(t)}|z^{(1)}, \dots, z^{(t)}).$$

The conditioning step conditions the distribution on an observation  $z^{(t+1)}$  using Bayes rule:

$$P(\sigma^{(t+1)}|z^{(1)}, \dots, z^{(t+1)}) \propto P(z^{(t+1)}|\sigma^{(t+1)})P(\sigma^{(t+1)}|z^{(1)}, \dots, z^{(t)}).$$

Since there are  $n!$  permutations, a single iteration of the algorithm requires  $O((n!)^2)$  flops and is consequently intractable for all but very small  $n$ . The approach that we advocate is to maintain a compact approximation to the true distribution based on the Fourier transform. As we discuss later, the Fourier based approximation is equivalent to maintaining a set of low-order marginals, rather than the full joint, which we regard as being analogous to an *Assumed Density Filter* (Boyer and Koller, 1998).

Although we use hidden Markov models and filtering as a running example, the approach we describe is useful for other probabilistic inference tasks over permutations, such as ranking objects and modeling user preferences. For example, operations such as marginalization and conditioning are fundamental and are widely applicable. In particular, conditioning using Bayes rule, one of the main topics of our paper, is one of the most fundamental probabilistic operations, and we provide a completely general formulation.

### 3. Probability Distributions over the Symmetric Group

A permutation on  $n$  elements is a one-to-one mapping of the set  $\{1, \dots, n\}$  into itself and can be written as a tuple,

$$\sigma = [\sigma(1) \ \sigma(2) \ \dots \ \sigma(n)],$$

where  $\sigma(i)$  denotes where the  $i$ th element is mapped under the permutation (called *one line notation*). For example,  $\sigma = [2 \ 3 \ 1 \ 4 \ 5]$  means that  $\sigma(1) = 2$ ,  $\sigma(2) = 3$ ,  $\sigma(3) = 1$ ,  $\sigma(4) = 4$ , and  $\sigma(5) = 5$ . The set of all permutations on  $n$  elements forms a group under the<sup>2</sup> operation of function composition—that is, if  $\sigma_1$  and  $\sigma_2$  are permutations, then

$$\sigma_1\sigma_2 = [\sigma_1(\sigma_2(1)) \ \sigma_1(\sigma_2(2)) \ \dots \ \sigma_1(\sigma_2(n))]$$

is itself a permutation. The set of all  $n!$  permutations is called the *symmetric group*, or just  $S_n$ .

We will actually notate the elements of  $S_n$  using the more standard *cycle notation*, in which a *cycle*  $(i, j, k, \dots, \ell)$  refers to the permutation which maps  $i$  to  $j$ ,  $j$  to  $k$ ,  $\dots$ , and finally  $\ell$  to  $i$ . Though not every permutation can be written as a single cycle, any permutation can always be written as a product of disjoint cycles. For example, the permutation  $\sigma = [2 \ 3 \ 1 \ 4 \ 5]$  written in cycle notation is  $\sigma = (1, 2, 3)(4)(5)$ . The number of elements in a cycle is called the *cycle length* and we typically drop the length 1 cycles in cycle notation when it creates no ambiguity—in our

2. See Appendix A for a list of the basic group theoretic definitions used in this paper.

example,  $\sigma = (1,2,3)(4)(5) = (1,2,3)$ . We refer to the identity permutation (which maps every element to itself) as  $\epsilon$ .

A probability distribution over permutations can be thought of as a joint distribution on the  $n$  random variables  $(\sigma(1), \dots, \sigma(n))$  subject to the *mutual exclusivity constraints* that  $P(\sigma : \sigma(i) = \sigma(j)) = 0$  whenever  $i \neq j$ . For example, in the identity management problem, Alice and Bob cannot both be in Track 1 simultaneously. Due to the fact that all of the  $\sigma(i)$  are coupled in the joint distribution, graphical models, which might have otherwise exploited an underlying conditional independence structure, are ineffective. Instead, our Fourier based approximation achieves compactness by exploiting the *algebraic structure* of the problem.

### 3.1 Compact Summary Statistics

While continuous distributions like Gaussians are typically summarized using moments (like mean and variance), or more generally, expected features, it is not immediately obvious how one might, for example, compute the ‘mean’ of a distribution over permutations. There is a simple method that might spring to mind, however, which is to think of the permutations as *permutation matrices* and to average the matrices instead.

**Example 1** For example, consider the two permutations  $\epsilon, (1,2) \in S_3$  ( $\epsilon$  is the identity and  $(1,2)$  swaps 1 and 2). We can associate the identity permutation  $\epsilon$  with the  $3 \times 3$  identity matrix, and similarly, we can associate the permutation  $(1,2)$  with the matrix:

$$(1,2) \mapsto \begin{bmatrix} 0 & 1 & 0 \\ 1 & 0 & 0 \\ 0 & 0 & 1 \end{bmatrix}.$$

The ‘average’ of  $\epsilon$  and  $(1,2)$  is therefore:

$$\frac{1}{2} \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} + \frac{1}{2} \begin{bmatrix} 0 & 1 & 0 \\ 1 & 0 & 0 \\ 0 & 0 & 1 \end{bmatrix} = \begin{bmatrix} 1/2 & 1/2 & 0 \\ 1/2 & 1/2 & 0 \\ 0 & 0 & 1 \end{bmatrix}.$$

As we will later show, computing the ‘mean’ (as described above) of a distribution over permutations,  $P$ , compactly summarizes  $P$  by storing a marginal distribution over each of  $\sigma(1), \sigma(2), \dots, \sigma(n)$ , which requires storing only  $O(n^2)$  numbers rather than the full  $O(n!)$  for the exact distribution. As an example, one possible summary might look like:

$$\hat{P} = \left[ \begin{array}{c|ccc} & \text{Alice} & \text{Bob} & \text{Cathy} \\ \hline \text{Track 1} & 2/3 & 1/6 & 1/6 \\ \text{Track 2} & 1/3 & 1/3 & 1/3 \\ \text{Track 3} & 0 & 1/2 & 1/2 \end{array} \right].$$

Such doubly stochastic “first-order summaries” have been studied in various settings (Shin et al., 2003; Helmbold and Warmuth, 2007). In identity management (Shin et al., 2003),<sup>3</sup> first-order sum-

3. Strictly speaking, a map from identities to tracks is not a permutation since a permutation always maps a set into itself. In fact, the set of all such identity-to-track assignments does not actually form a group since there is no way to compose any two such assignments to obtain a legitimate group operation. We abuse the notation by referring to these assignments as a group, but really the elements of the group here should be thought of as the ‘deviation’ from the original identity-to-track assignment (where only the tracks are permuted, for example, when they are confused). In the group theoretic language, there is a faithful group action of  $S_n$  on the set of all identity-to-track assignments.

maries maintain, for example,

$$\begin{aligned} P(\text{Alice is at Track 1}) &= 2/3, \\ P(\text{Bob is at Track 3}) &= 1/2. \end{aligned}$$

What cannot be captured by first-order summaries however, are the higher order statements like:

$$P(\text{Alice is in Track 1 and Bob is in Track 2}) = 0.$$

Over the next two sections, we will show that the first-order summary of a distribution  $P(\sigma)$  can equivalently be viewed as the lowest frequency coefficients of the Fourier transform of  $P(\sigma)$ , and that by considering higher frequencies, we can capture higher order marginal probabilities in a principled fashion. Furthermore, the Fourier theoretic perspective, as we will show, provides a natural framework for formulating inference operations with respect to our compact summaries. In a nutshell, we will view the prediction/rollup step as a convolution and the conditioning step as a pointwise product—then we will formulate the two inference operations in the Fourier domain as a pointwise product and convolution, respectively.

## 4. The Fourier Transform on Finite Groups

Over the last fifty years, the Fourier Transform has been ubiquitously applied to everything digital, particularly with the invention of the Fast Fourier Transform (Cooley and Tukey, 1965; Rockmore, 2000). On the real line, the Fourier Transform is a well-studied method for decomposing a function into a sum of sine and cosine terms over a spectrum of frequencies. Perhaps less familiar to the machine learning community though, is its group theoretic generalization. In this section we review group theoretic generalizations of the Fourier transform with an eye towards approximating functions on  $S_n$ . None of the results stated in this section or the next are original. Noncommutative generalizations of the Fourier transform have been studied quite extensively throughout the last century from both the mathematics (Lang, 1965) and physics communities (Chen, 1989). Applications to permutations were first pioneered by Persi Diaconis who studied problems in card shuffling and since then, there have been many papers on related topics in probability and statistics. For further information, see Diaconis (1988) and Terras (1999).

### 4.1 Group Representation Theory

The generalized definition of the Fourier Transform relies on the theory of group representations, which formalize the concept of associating permutations with matrices and are used to construct a complete basis for the space of functions on a group  $G$ , thus also playing a role analogous to that of sinusoids on the real line.

**Definition 1** A representation of a group  $G$  is a map  $\rho$  from  $G$  to a set of invertible  $d_\rho \times d_\rho$  (complex) matrix operators ( $\rho : G \rightarrow \mathbb{C}^{d_\rho \times d_\rho}$ ) which preserves algebraic structure in the sense that for all  $\sigma_1, \sigma_2 \in G$ ,  $\rho(\sigma_1 \sigma_2) = \rho(\sigma_1) \cdot \rho(\sigma_2)$ . The matrices which lie in the image of  $\rho$  are called the representation matrices, and we will refer to  $d_\rho$  as the degree of the representation.

The requirement that  $\rho(\sigma_1\sigma_2) = \rho(\sigma_1) \cdot \rho(\sigma_2)$  is analogous to the property that  $e^{i(\theta_1+\theta_2)} = e^{i\theta_1} \cdot e^{i\theta_2}$  for the conventional sinusoidal basis. Each matrix entry,  $\rho_{ij}(\sigma)$  defines some function over  $S_n$ :

$$\rho(\sigma) = \begin{bmatrix} \rho_{11}(\sigma) & \rho_{12}(\sigma) & \cdots & \rho_{1d_\rho}(\sigma) \\ \rho_{21}(\sigma) & \rho_{22}(\sigma) & \cdots & \rho_{2d_\rho}(\sigma) \\ \vdots & \vdots & \ddots & \vdots \\ \rho_{d_\rho 1}(\sigma) & \rho_{d_\rho 2}(\sigma) & \cdots & \rho_{d_\rho d_\rho}(\sigma) \end{bmatrix},$$

and consequently, each representation  $\rho$  simultaneously defines a set of  $d_\rho^2$  functions over  $S_n$ . We will eventually think of group representations as the set of Fourier basis functions onto which we can project arbitrary functions.

Before moving onto examples, we make several remarks about the generality of this paper. First, while our paper is primarily focused on the symmetric group, many of its results hold for arbitrary finite groups. For example, there are a variety of finite groups that have been studied in applications, like metacyclic groups (Willsky, 1978), wreath product groups (Foote et al., 2004), etc. However, while some of these results will even extend with minimal effort to more general cases, such as locally compact groups, the assumption in all of the following results will be that  $G$  is finite, even if it is not explicitly stated. Specifically, most of the results in Sections 4, 6, and Appendix D.2 are intended to hold over any finite group, while the results of the remaining sections are specific to probabilistic inference over the symmetric group. Secondly, given an arbitrary finite group  $G$ , some of the algebraic results that we use require that the underlying field be the complex numbers. For the particular case of the symmetric group, however, we can in fact assume that the representations are real-valued matrices. Thus, throughout the paper, we will explicitly assume that the representations are real-valued.<sup>4</sup>

**Example 2** *We begin by showing three examples of representations on the symmetric group.*

1. *The simplest example of a representation is called the trivial representation  $\rho_{(n)} : S_n \rightarrow \mathbb{R}^{1 \times 1}$ , which maps each element of the symmetric group to 1, the multiplicative identity on the real numbers. The trivial representation is actually defined for every group, and while it may seem unworthy of mention, it plays the role of the constant basis function in the Fourier theory.*
2. *The first-order permutation representation of  $S_n$ , which we alluded to in Example 1, is the degree  $n$  representation,  $\tau_{(n-1,1)}$  (we explain the terminology in Section 5), which maps a permutation  $\sigma$  to its corresponding permutation matrix given by  $[\tau_{(n-1,1)}(\sigma)]_{ij} = \mathbf{1}\{\sigma(j) = i\}$ . For example, the first-order permutation representation on  $S_3$  is given by:*

$$\begin{aligned} \tau_{(2,1)}(\epsilon) &= \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} & \tau_{(2,1)}(1,2) &= \begin{bmatrix} 0 & 1 & 0 \\ 1 & 0 & 0 \\ 0 & 0 & 1 \end{bmatrix} & \tau_{(2,1)}(2,3) &= \begin{bmatrix} 1 & 0 & 0 \\ 0 & 0 & 1 \\ 0 & 1 & 0 \end{bmatrix} \\ \tau_{(2,1)}(1,3) &= \begin{bmatrix} 0 & 0 & 1 \\ 0 & 1 & 0 \\ 1 & 0 & 0 \end{bmatrix} & \tau_{(2,1)}(1,2,3) &= \begin{bmatrix} 0 & 0 & 1 \\ 1 & 0 & 0 \\ 0 & 1 & 0 \end{bmatrix} & \tau_{(2,1)}(1,3,2) &= \begin{bmatrix} 0 & 1 & 0 \\ 0 & 0 & 1 \\ 1 & 0 & 0 \end{bmatrix} \end{aligned}$$

---

4. To recover similar results for more complex-valued representations, one would have to replace matrix transposes by adjoints, etc.



3. The alternating representation of  $S_n$ , maps a permutation  $\sigma$  to the determinant of  $\tau_{(n-1,1)}(\sigma)$ , which is  $+1$  if  $\sigma$  can be equivalently written as the composition of an even number of pairwise swaps, and  $-1$  otherwise. We write the alternating representation as  $\rho_{(1,\dots,1)}$  with  $n$  1's in the subscript. For example, on  $S_4$ , we have:

$$\rho_{(1,1,1,1)}((1,2,3)) = \rho_{(1,1,1,1)}((13)(12)) = +1.$$

The alternating representation can be interpreted as the ‘highest frequency’ basis function on the symmetric group, intuitively due to its high sensitivity to swaps. For example, if  $\tau_{(1,\dots,1)}(\sigma) = 1$ , then  $\tau_{(1,\dots,1)}((12)\sigma) = -1$ . In identity management, it may be reasonable to believe that the joint probability over all  $n$  identity labels should only change by a little if just two objects are mislabeled due to swapping—in this case, ignoring the basis function corresponding to the alternating representation should still provide an accurate approximation to the joint distribution.

In general, a representation corresponds to an overcomplete set of functions and therefore does not constitute a valid basis for any subspace of functions. For example, the set of nine functions on  $S_3$  corresponding to  $\tau_{(2,1)}$  span only four dimensions, because there are six normalization constraints (three on the row sums and three on the column sums), of which five are independent—and so there are five redundant dimensions. To find a valid complete basis for the space of functions on  $S_n$ , we will need to find a family of representations whose basis functions are independent, and span the entire  $n!$ -dimensional space of functions.

In the following two definitions, we will provide two methods for constructing a new representation from old ones such that the set of functions on  $S_n$  corresponding to the new representation is linearly *dependent* on the old representations. Somewhat surprisingly, it can be shown that dependencies which arise amongst the representations can always be recognized in a certain sense, to come from the two possible following sources (Serre, 1977).

## Definition 2

1. **Equivalence.** Given a representation  $\rho_1$  and an invertible matrix  $C$ , one can define a new representation  $\rho_2$  by “changing the basis” for  $\rho_1$ :

$$\rho_2(\sigma) \triangleq C^{-1} \cdot \rho_1(\sigma) \cdot C. \quad (1)$$

We say, in this case, that  $\rho_1$  and  $\rho_2$  are equivalent as representations (written  $\rho_1 \equiv \rho_2$ , as opposed to  $\rho_1 = \rho_2$ ), and the matrix  $C$  is known as the intertwining operator. Note that  $d_{\rho_1} = d_{\rho_2}$ .

It can be checked that the functions corresponding to  $\rho_2$  can be reconstructed from those corresponding to  $\rho_1$ . For example, if  $C$  is a permutation matrix, the matrix entries of  $\rho_2$  are exactly the same as the matrix entries of  $\rho_1$ , only permuted.

2. **Direct Sum.** Given two representations  $\rho_1$  and  $\rho_2$ , we can always form a new representation, which we will write as  $\rho_1 \oplus \rho_2$ , by defining:

$$\rho_1 \oplus \rho_2(\sigma) \triangleq \left[ \begin{array}{c|c} \rho_1(\sigma) & 0 \\ \hline 0 & \rho_2(\sigma) \end{array} \right].$$

$\rho_1 \oplus \rho_2$  is called the direct sum representation. For example, the direct sum of two copies of the trivial representation is:

$$\rho_{(n)} \oplus \rho_{(n)}(\sigma) = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix},$$

with four corresponding functions on  $S_n$ , each of which is clearly dependent upon the trivial representation itself.

Most representations can be seen as being equivalent to a direct sum of strictly smaller representations. Whenever a representation  $\rho$  can be decomposed as  $\rho \equiv \rho_1 \oplus \rho_2$ , we say that  $\rho$  is *reducible*. As an example, we now show that the first-order permutation representation is a reducible representation.

**Example 3** Instead of using the standard basis vectors  $\{e_1, e_2, e_3\}$ , the first-order permutation representation for  $S_3$ ,  $\tau_{(2,1)} : S_3 \rightarrow \mathbb{C}^{3 \times 3}$ , can be equivalently written with respect to a new basis  $\{v_1, v_2, v_3\}$ , where:

$$\begin{aligned} v_1 &= \frac{e_1 + e_2 + e_3}{|e_1 + e_2 + e_3|}, \\ v_2 &= \frac{-e_1 + e_2}{|-e_1 + e_2|}, \\ v_3 &= \frac{-e_1 - e_2 + 2e_3}{|-e_1 - e_2 + 2e_3|}. \end{aligned}$$

To ‘change the basis’, we write the new basis vectors as columns in a matrix  $C$ :

$$C = \begin{bmatrix} | & | & | \\ v_1 & v_2 & v_3 \\ | & | & | \end{bmatrix} = \begin{bmatrix} \frac{1}{\sqrt{3}} & -\frac{\sqrt{2}}{2} & -\frac{1}{\sqrt{6}} \\ \frac{1}{\sqrt{3}} & \frac{\sqrt{2}}{2} & -\frac{1}{\sqrt{6}} \\ \frac{1}{\sqrt{3}} & 0 & \frac{2}{\sqrt{6}} \end{bmatrix},$$

and conjugate the representation  $\tau_{(2,1)}$  by  $C$  (as in Equation 1) to obtain the equivalent representation  $C^{-1} \cdot \tau_{(2,1)}(\sigma) \cdot C$ :

$$\begin{aligned} C^{-1} \cdot \tau_{(2,1)}(\epsilon) \cdot C &= \left[ \begin{array}{c|cc} 1 & 0 & 0 \\ \hline 0 & 1 & 0 \\ 0 & 0 & 1 \end{array} \right] & C^{-1} \cdot \tau_{(2,1)}(1,2) \cdot C &= \left[ \begin{array}{c|cc} 1 & 0 & 0 \\ \hline 0 & -1 & 0 \\ 0 & 0 & 1 \end{array} \right] \\ C^{-1} \cdot \tau_{(2,1)}(2,3) \cdot C &= \left[ \begin{array}{c|cc} 1 & 0 & 0 \\ \hline 0 & 1/2 & \sqrt{3}/2 \\ 0 & \sqrt{3}/2 & -1/2 \end{array} \right] & C^{-1} \cdot \tau_{(2,1)}(1,3) \cdot C &= \left[ \begin{array}{c|cc} 1 & 0 & 0 \\ \hline 0 & 1/2 & -\sqrt{3}/2 \\ 0 & -\sqrt{3}/2 & -1/2 \end{array} \right] \\ C^{-1} \cdot \tau_{(2,1)}(1,2,3) \cdot C &= \left[ \begin{array}{c|cc} 1 & 0 & 0 \\ \hline 0 & -1/2 & -\sqrt{3}/2 \\ 0 & \sqrt{3}/2 & -1/2 \end{array} \right] & C^{-1} \cdot \tau_{(2,1)}(1,3,2) \cdot C &= \left[ \begin{array}{c|cc} 1 & 0 & 0 \\ \hline 0 & -1/2 & \sqrt{3}/2 \\ 0 & -\sqrt{3}/2 & -1/2 \end{array} \right] \end{aligned}$$

The interesting property of this particular basis is that the new representation matrices all appear to be the direct sum of two smaller representations, a trivial representation,  $\rho_{(3)}$  as the top left block, and a degree 2 representation in the bottom right which we will refer to as  $\rho_{(2,1)}$ .

Geometrically, the representation  $\rho_{(2,1)}$  can also be thought of as the group of rigid symmetries of the equilateral triangle with vertices:

$$P_1 = \begin{bmatrix} \sqrt{3}/2 \\ 1/2 \end{bmatrix}, P_2 = \begin{bmatrix} -\sqrt{3}/2 \\ 1/2 \end{bmatrix}, P_3 = \begin{bmatrix} 0 \\ -1 \end{bmatrix}.$$

The matrix  $\rho_{(2,1)}(1, 2)$  acts on the triangle by reflecting about the  $x$ -axis, and  $\rho_{(2,1)}(1, 2, 3)$  by a  $\pi/3$  counter-clockwise rotation.

In general, there are infinitely many reducible representations. For example, given any dimension  $d$ , there is a representation which maps every element of a group  $G$  to the  $d \times d$  identity matrix (the direct sum of  $d$  copies of the trivial representation). However, for any finite group, there exists a finite collection of atomic representations which can be used to build up any other representation (up to equivalence) using the direct sum operation. These representations are referred to as the *irreducibles* of a group, and they are defined simply to be the collection of representations (up to equivalence) which are not reducible. It can be shown that any (complex) representation of a finite group  $G$  is equivalent to a direct sum of irreducibles (Diaconis, 1988), and hence, for any representation  $\tau$ , there exists a matrix  $C$  for which

$$C^{-1} \cdot \tau \cdot C = \bigoplus_{\rho} \bigoplus_{j=1}^{z_{\rho}} \rho, \quad (2)$$

where  $\rho$  ranges over all distinct irreducible representations of the group  $G$ , and the inner  $\oplus$  refers to some finite number ( $z_{\rho}$ ) of copies of each irreducible  $\rho$ .

As it happens, there are only three irreducible representations of  $S_3$  (Diaconis, 1988), up to equivalence: the trivial representation  $\rho_{(3)}$ , the degree 2 representation  $\rho_{(2,1)}$ , and the alternating representation  $\rho_{(1,1,1)}$ . The complete set of irreducible representation matrices of  $S_3$  are shown in the Table 2. Unfortunately, the analysis of the irreducible representations for  $n > 3$  is far more complicated and we postpone this more general discussion for Section 5.

## 4.2 The Fourier Transform

The link between group representation theory and Fourier analysis is given by the celebrated *Peter-Weyl theorem* (Diaconis, 1988; Terras, 1999; Sagan, 2001) which says that the matrix entries of the irreducibles of  $G$  form a *complete* set of *orthogonal* basis functions on  $G$ .<sup>5</sup> The space of functions on  $S_3$ , for example, is orthogonally spanned by the  $3!$  functions  $\rho_{(3)}(\sigma)$ ,  $[\rho_{(2,1)}(\sigma)]_{1,1}$ ,  $[\rho_{(2,1)}(\sigma)]_{1,2}$ ,  $[\rho_{(2,1)}(\sigma)]_{2,1}$ ,  $[\rho_{(2,1)}(\sigma)]_{2,2}$  and  $\rho_{(1,1,1)}(\sigma)$ , where  $[\rho(\sigma)]_{ij}$  denotes the  $(i, j)$  entry of the matrix  $\rho(\sigma)$ .

As a replacement for projecting a function  $f$  onto a complete set of sinusoidal basis functions (as one would do on the real line), the Peter-Weyl theorem suggests instead to project onto the basis provided by the irreducibles of  $G$ . As on the real line, this projection can be done by computing the inner product of  $f$  with each element of the basis, and we define this operation to be the generalized form of the Fourier Transform.

5. Technically the Peter-Weyl result, as stated here, is only true if all of the representation matrices are unitary. That is,  $\rho(\sigma)^* \rho(\sigma) = I$  for all  $\sigma \in S_n$ , where the matrix  $A^*$  is the conjugate transpose of  $A$ . For the case of real-valued (as opposed to complex-valued) matrices, however, the definitions of unitary and orthogonal matrices coincide.

While most representations are not unitary, there is a standard result from representation theory which shows that for *any* representation of  $G$ , there exists an equivalent unitary representation.

$\sigma$	$\rho_{(3)}$	$\rho_{(2,1)}$	$\rho_{(1,1,1)}$
$\epsilon$	1	$\begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$	1
(1,2)	1	$\begin{bmatrix} -1 & 0 \\ 0 & 1 \end{bmatrix}$	-1
(2,3)	1	$\begin{bmatrix} 1/2 & \sqrt{3}/2 \\ \sqrt{3}/2 & -1/2 \end{bmatrix}$	-1
(1,3)	1	$\begin{bmatrix} 1/2 & -\sqrt{3}/2 \\ -\sqrt{3}/2 & -1/2 \end{bmatrix}$	-1
(1,2,3)	1	$\begin{bmatrix} -1/2 & -\sqrt{3}/2 \\ \sqrt{3}/2 & -1/2 \end{bmatrix}$	1
(1,3,2)	1	$\begin{bmatrix} -1/2 & \sqrt{3}/2 \\ -\sqrt{3}/2 & -1/2 \end{bmatrix}$	1

Table 2: The irreducible representation matrices of  $S_3$ .

**Definition 3** Let  $f : G \rightarrow \mathbb{R}$  be any function on a group  $G$  and let  $\rho$  be any representation on  $G$ . The Fourier Transform of  $f$  at the representation  $\rho$  is defined to be the matrix of coefficients:

$$\hat{f}_\rho = \sum_{\sigma} f(\sigma)\rho(\sigma).$$

The collection of Fourier Transforms at all irreducible representations of  $G$  form the Fourier Transform of  $f$ .

There are two important points which distinguish this Fourier Transform from its familiar formulation on the real line—first, the outputs of the transform are matrix-valued, and second, the inputs to  $\hat{f}$  are *representations* of  $G$  rather than real numbers. As in the familiar formulation, the Fourier Transform is invertible and the inversion formula is explicitly given by the Fourier Inversion Theorem.

**Theorem 4 (Fourier Inversion Theorem)**

$$f(\sigma) = \frac{1}{|G|} \sum_{\lambda} d_{\rho_\lambda} \text{Tr} \left[ \hat{f}_{\rho_\lambda}^T \cdot \rho_\lambda(\sigma) \right], \tag{3}$$

where  $\lambda$  indexes over the collection of irreducibles of  $G$ .

Note that the trace term in the inverse Fourier Transform is just the ‘matrix dot product’ between  $\hat{f}_{\rho_\lambda}$  and  $\rho_\lambda(\sigma)$ , since  $\text{Tr} [A^T \cdot B] = \langle \text{vec}(A), \text{vec}(B) \rangle$ , where by  $\text{vec}$  we mean mapping a matrix to a vector on the same elements arranged in column-major order.

We now provide several examples for intuition. For functions on the real line, the Fourier Transform at zero frequency gives the DC component of a signal. The same holds true for functions on a group; If  $f : G \rightarrow \mathbb{R}$  is any function, then since  $\rho_{(n)} = 1$ , the Fourier Transform of  $f$  at the

trivial representation is constant, with  $\hat{f}_{\rho(n)} = \sum_{\sigma} f(\sigma)$ . Thus, for any probability distribution  $P$ , we have  $\hat{P}_{\rho(n)} = 1$ . If  $P$  were the uniform distribution, then  $\hat{P}_{\rho} = 0$  at every irreducible  $\rho$  except at the trivial representation.

The Fourier Transform at  $\tau_{(n-1,1)}$  also has a simple interpretation:

$$[\hat{f}_{\tau_{(n-1,1)}}]_{ij} = \sum_{\sigma \in S_n} f(\sigma) [\tau_{(n-1,1)}(\sigma)]_{ij} = \sum_{\sigma \in S_n} f(\sigma) \mathbb{1}\{\sigma(j) = i\} = \sum_{\sigma: \sigma(j)=i} f(\sigma).$$

The set  $\Delta_{ij} = \{\sigma : \sigma(j) = i\}$  is the set of the  $(n-1)!$  possible permutations which map element  $j$  to  $i$ . In identity management,  $\Delta_{ij}$  can be thought of as the set of assignments which, for example, have Alice at Track 1. If  $P$  is a distribution, then  $\hat{P}_{\tau_{(n-1,1)}}$  is a matrix of *first-order* marginal probabilities, where the  $(i, j)$ -th element is the marginal probability that a random permutation drawn from  $P$  maps element  $j$  to  $i$ .

**Example 4** Consider the following probability distribution on  $S_3$ :

$\sigma$	$\epsilon$	(1,2)	(2,3)	(1,3)	(1,2,3)	(1,3,2)
$P(\sigma)$	1/3	1/6	1/3	0	1/6	0

The set of all first order marginal probabilities is given by the Fourier transform at  $\tau_{(2,1)}$ :

$$\hat{P}_{\tau_{(2,1)}} = \begin{bmatrix} & \mathbf{A} & \mathbf{B} & \mathbf{C} \\ \mathbf{1} & 2/3 & 1/6 & 1/6 \\ \mathbf{2} & 1/3 & 1/3 & 1/3 \\ \mathbf{3} & 0 & 1/2 & 1/2 \end{bmatrix}.$$

In the above matrix, each column  $j$  represents a marginal distribution over the possible tracks that identity  $j$  can map to under a random draw from  $P$ . We see, for example, that Alice is at Track 1 with probability 2/3, or at Track 2 with probability 1/3. Simultaneously, each row  $i$  represents a marginal distribution over the possible identities that could have been mapped to track  $i$  under a random draw from  $P$ . In our example, Bob and Cathy are equally likely to be in Track 3, but Alice is definitely not in Track 3. Since each row and each column is itself a distribution, the matrix  $\hat{P}_{\tau_{(2,1)}}$  must be doubly stochastic. We will elaborate on the consequences of this observation later.

The Fourier transform of the same distribution at all irreducibles is:

$$\hat{P}_{\rho(3)} = 1, \quad \hat{P}_{\rho(2,1)} = \begin{bmatrix} 1/4 & \sqrt{3}/4 \\ \sqrt{3}/4 & 1/4 \end{bmatrix}, \quad \hat{P}_{\rho(1,1,1)} = 0.$$

The first-order permutation representation,  $\tau_{(n-1,1)}$ , captures the statistics of how a random permutation acts on a single object irrespective of where all of the other  $n-1$  objects are mapped, and in doing so, compactly summarizes the distribution with only  $O(n^2)$  numbers. Unfortunately, as mentioned in Section 3, the Fourier transform at the first-order permutation representation cannot capture more complicated statements like:

$$P(\text{Alice and Bob occupy Tracks 1 and 2}) = 0.$$

To avoid collapsing away so much information, we might define richer summary statistics that might capture ‘higher-order’ effects. We define the *second-order unordered permutation representation* by:

$$[\tau_{(n-2,2)}(\sigma)]_{\{i,j\},\{k,\ell\}} = \mathbb{1} \{ \sigma(\{k, \ell\}) = \{i, j\} \},$$

where we index the matrix rows and columns by unordered pairs  $\{i, j\}$ . The condition inside the indicator function states that the representation captures whether the pair of objects  $\{k, \ell\}$  maps to the pair  $\{i, j\}$ , but is indifferent with respect to the ordering; that is, either  $k \mapsto i$  and  $\ell \mapsto j$ , or,  $k \mapsto j$  and  $\ell \mapsto i$ .

**Example 5** For  $n = 4$ , there are six possible unordered pairs:  $\{1, 2\}, \{1, 3\}, \{1, 4\}, \{2, 3\}, \{2, 4\}$ , and  $\{3, 4\}$ . The matrix representation of the permutation  $(1, 2, 3)$  is:

$$\tau_{(2,2)}(1, 2, 3) = \begin{array}{c|cccccc} & \{1, 2\} & \{1, 3\} & \{1, 4\} & \{2, 3\} & \{2, 4\} & \{3, 4\} \\ \hline \{1, 2\} & 0 & 0 & 0 & 1 & 0 & 0 \\ \{1, 3\} & 1 & 0 & 0 & 0 & 0 & 0 \\ \{1, 4\} & 0 & 0 & 0 & 0 & 1 & 0 \\ \{2, 3\} & 0 & 1 & 0 & 0 & 0 & 0 \\ \{2, 4\} & 0 & 0 & 0 & 0 & 0 & 1 \\ \{3, 4\} & 0 & 0 & 1 & 0 & 0 & 0 \end{array}.$$

The *second order ordered permutation representation*,  $\tau_{(n-2,1,1)}$ , is defined similarly:

$$[\tau_{(n-2,1,1)}(\sigma)]_{(i,j),(k,\ell)} = \mathbb{1} \{ \sigma((k, \ell)) = (i, j) \},$$

where  $(k, \ell)$  denotes an *ordered* pair. Therefore,  $[\tau_{(n-2,1,1)}(\sigma)]_{(i,j),(k,\ell)}$  is 1 if and only if  $\sigma$  maps  $k$  to  $i$  and  $\ell$  to  $j$ .

As in the first-order case, the Fourier transform of a probability distribution at  $\tau_{(n-2,2)}$ , returns a matrix of marginal probabilities of the form:  $P(\sigma : \sigma(\{k, \ell\}) = \{i, j\})$ , which captures statements like, "Alice and Bob occupy Tracks 1 and 2 with probability 1/2". Similarly, the Fourier transform at  $\tau_{(n-2,1,1)}$  returns a matrix of marginal probabilities of the form  $P(\sigma : \sigma((k, \ell)) = (i, j))$ , which captures statements like, "Alice is in Track 1 and Bob is in Track 2 with probability 9/10".

We can go further and define third-order representations, fourth-order representations, and so on. In general however, the permutation representations as they have been defined above are reducible, intuitively due to the fact that it is possible to recover lower order marginal probabilities from higher order marginal probabilities. For example, one can recover the normalization constant (corresponding to the trivial representation) from the first order matrix of marginals by summing across either the rows or columns, and the first order marginal probabilities from the second order marginal probabilities by summing across appropriate matrix entries. To truly leverage the machinery of Fourier analysis, it is important to understand the Fourier transform at the irreducibles of the symmetric group, and in the next section, we show how to derive the irreducible representations of the Symmetric group by first defining permutation representations, then “subtracting off the lower-order effects”.

## 5. Representation Theory on the Symmetric Group

In this section, we provide a brief introduction to the representation theory of the Symmetric group. Rather than giving a fully rigorous treatment of the subject, our goal is to give some intuition about

the kind of information which can be captured by the irreducible representations of  $S_n$ . Roughly speaking, we will show that Fourier transforms on the Symmetric group, instead of being indexed by frequencies, are indexed by *partitions* of  $n$  (tuples of numbers which sum to  $n$ ), and certain partitions correspond to more complex basis functions than others. For proofs, we point the reader to consult: Diaconis (1989), James and Kerber (1981), Sagan (2001) and Vershik and Okounkov (2006).

Instead of the singleton or pairwise marginals which were described in the previous section, we will now focus on using the Fourier coefficients of a distribution to query a much wider class of marginal probabilities. As an example, we will be able to compute the following (more complicated) marginal probability on  $S_6$  using Fourier coefficients:

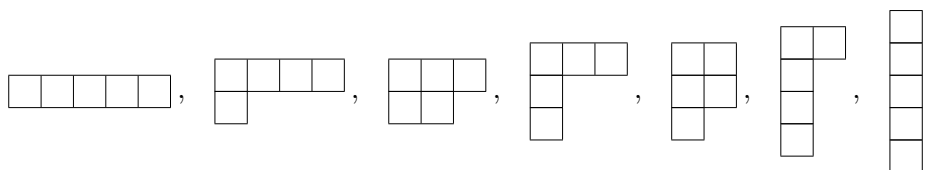
$$P\left(\sigma : \sigma\left(\left\{\begin{array}{|c|c|c|} \hline 1 & 2 & 3 \\ \hline 4 & 5 & \\ \hline 6 & & \end{array}\right\}\right) = \left\{\begin{array}{|c|c|c|} \hline 1 & 2 & 6 \\ \hline 4 & 5 & \\ \hline 3 & & \end{array}\right\}, \tag{4}$$

which we interpret as the joint marginal probability that the rows of the diagram on the left map to corresponding rows on the right as unordered sets. In other words, Equation 4 is the joint probability that *unordered* set  $\{1,2,3\}$  maps to  $\{1,2,6\}$ , the unordered pair  $\{4,5\}$  maps to  $\{4,5\}$ , and the singleton  $\{6\}$  maps to  $\{3\}$ .

The diagrams in Equation 4 are known as *Ferrer's diagrams* and are commonly used to visualize *partitions* of  $n$ , which are defined to be unordered tuples of positive integers,  $\lambda = (\lambda_1, \dots, \lambda_\ell)$ , which sum to  $n$ . For example,  $\lambda = (3, 2)$  is a partition of  $n = 5$  since  $3 + 2 = 5$ . Usually we write partitions as weakly decreasing sequences by convention, so the partitions of  $n = 5$  are:

$$(5), (4, 1), (3, 2), (3, 1, 1), (2, 2, 1), (2, 1, 1, 1), (1, 1, 1, 1, 1),$$

and their respective Ferrers diagrams are:



A *Young tabloid* is an assignment of the numbers  $\{1, \dots, n\}$  to the boxes of a Ferrers diagram for a partition  $\lambda$ , where each row represents an unordered set. There are 6 Young tabloids corresponding to the partition  $\lambda = (2, 2)$ , for example:

$$\left\{\begin{array}{|c|c|} \hline 1 & 2 \\ \hline 3 & 4 \\ \hline \end{array}\right\}, \left\{\begin{array}{|c|c|} \hline 1 & 3 \\ \hline 2 & 4 \\ \hline \end{array}\right\}, \left\{\begin{array}{|c|c|} \hline 1 & 4 \\ \hline 2 & 3 \\ \hline \end{array}\right\}, \left\{\begin{array}{|c|c|} \hline 2 & 3 \\ \hline 1 & 4 \\ \hline \end{array}\right\}, \left\{\begin{array}{|c|c|} \hline 2 & 4 \\ \hline 1 & 3 \\ \hline \end{array}\right\}, \left\{\begin{array}{|c|c|} \hline 3 & 4 \\ \hline 1 & 2 \\ \hline \end{array}\right\}.$$

The Young tabloid,  $\begin{array}{|c|c|} \hline 1 & 2 \\ \hline 3 & 4 \\ \hline \end{array}$ , for example, represents the two unordered sets  $\{1,2\}$  and  $\{3,4\}$ , and if we were interested in computing the joint probability that  $\sigma(\{1,2\}) = \{3,4\}$  and  $\sigma(\{3,4\}) = \{1,2\}$ , then we could write the problem in terms of Young tabloids as:

$$P\left(\sigma : \sigma\left(\left\{\begin{array}{|c|c|} \hline 1 & 2 \\ \hline 3 & 4 \\ \hline \end{array}\right\}\right) = \left\{\begin{array}{|c|c|} \hline 3 & 4 \\ \hline 1 & 2 \\ \hline \end{array}\right\}\right).$$

In general, we will be able to use the Fourier coefficients at irreducible representations to compute the marginal probabilities of Young tabloids. As we shall see, with the help of the *James Submodule theorem* (James and Kerber, 1981), the marginals corresponding to “simple” partitions will require very few Fourier coefficients to compute, which is one of the main strengths of working in the Fourier domain.

**Example 6** *Imagine three separate rooms containing two tracks each, in which Alice and Bob are in room 1 occupying Tracks 1 and 2; Cathy and David are in room 2 occupying Tracks 3 and 4; and Eric and Frank are in room 3 occupying Tracks 5 and 6, but we are not able to distinguish which person is at which track in any of the rooms. Then*

$$P\left(\sigma : \left(\left\{\begin{array}{|c|c|} \hline A & B \\ \hline C & D \\ \hline E & F \\ \hline \end{array}\right\}\right) \rightarrow \left\{\begin{array}{|c|c|} \hline 1 & 2 \\ \hline 3 & 4 \\ \hline 5 & 6 \\ \hline \end{array}\right\}\right) = 1.$$

It is in fact, possible to recast the first-order marginals which were described in the previous section in the language of Young tabloids by noticing that, for example, if 1 maps to 1, then the unordered set  $\{2, \dots, n\}$  must map to  $\{2, \dots, n\}$  since permutations are one-to-one mappings. The marginal probability that  $\sigma(1) = 1$ , then, is equal to the marginal probability that  $\sigma(1) = 1$  and  $\sigma(\{2, \dots, n\}) = \{2, \dots, n\}$ . If  $n = 6$ , then the marginal probability written using Young tabloids is:

$$P\left(\sigma : \sigma\left(\left\{\begin{array}{|c|c|c|c|c|c|} \hline 2 & 3 & 4 & 5 & 6 \\ \hline 1 \\ \hline \end{array}\right\}\right) = \left\{\begin{array}{|c|c|c|c|c|c|} \hline 2 & 3 & 4 & 5 & 6 \\ \hline 1 \\ \hline \end{array}\right\}\right).$$

The first-order marginal probabilities correspond, therefore, to the marginal probabilities of Young tabloids of shape  $\lambda = (n - 1, 1)$ .

Likewise, the second-order unordered marginals correspond to Young tabloids of shape  $\lambda = (n - 2, 2)$ . If  $n = 6$  again, then the marginal probability that  $\{1, 2\}$  maps to  $\{2, 4\}$  corresponds to the following marginal probability for tabloids:

$$P\left(\sigma : \sigma\left(\left\{\begin{array}{|c|c|c|c|} \hline 3 & 4 & 5 & 6 \\ \hline 1 & 2 \\ \hline \end{array}\right\}\right) = \left\{\begin{array}{|c|c|c|c|} \hline 1 & 3 & 5 & 6 \\ \hline 2 & 4 \\ \hline \end{array}\right\}\right).$$

The second-order *ordered* marginals are captured at the partition  $\lambda = (n - 2, 1, 1)$ . For example, the marginal probability that  $\{1\}$  maps to  $\{2\}$  and  $\{2\}$  maps to  $\{4\}$  is given by:

$$P\left(\sigma : \sigma\left(\left\{\begin{array}{|c|c|c|c|} \hline 3 & 4 & 5 & 6 \\ \hline 1 \\ \hline 2 \\ \hline \end{array}\right\}\right) = \left\{\begin{array}{|c|c|c|c|} \hline 1 & 3 & 5 & 6 \\ \hline 2 \\ \hline 4 \\ \hline \end{array}\right\}\right).$$

And finally, we remark that the  $(1, \dots, 1)$  partition of  $n$  recovers all original probabilities since it asks for a joint distribution over  $\sigma(1), \dots, \sigma(n)$ . The corresponding matrix of marginals has  $n! \times n!$  entries (though there will only be  $n!$  distinct probabilities).

To see how the marginal probabilities of Young tabloids of shape  $\lambda$  can be thought of as Fourier coefficients, we will define a representation (which we call the *permutation representation*) associated with  $\lambda$  and show that the Fourier transform of a distribution at a permutation representation



gives marginal probabilities. We begin by fixing an ordering on the set of possible Young tabloids,  $\{t_1\}, \{t_2\}, \dots$ , and define the permutation representation  $\tau_\lambda(\sigma)$  to be the matrix:

$$[\tau_\lambda(\sigma)]_{ij} = \begin{cases} 1 & \text{if } \sigma(\{t_j\}) = \{t_i\} \\ 0 & \text{otherwise} \end{cases}.$$

It can be checked that the function  $\tau_\lambda$  is indeed a valid representation of the Symmetric group, and therefore we can compute Fourier coefficients at  $\tau_\lambda$ . If  $P(\sigma)$  is a probability distribution, then

$$\begin{aligned} [\widehat{P}_{\tau_\lambda}]_{ij} &= \sum_{\sigma \in \mathcal{S}_n} P(\sigma) [\tau_\lambda(\sigma)]_{ij}, \\ &= \sum_{\{\sigma: \sigma(\{t_j\}) = \{t_i\}\}} P(\sigma), \\ &= P(\sigma : \sigma(\{t_j\}) = \{t_i\}), \end{aligned}$$

and therefore, *the matrix of marginals corresponding to Young tabloids of shape  $\lambda$  is given exactly by the Fourier transform at the representation  $\tau_\lambda$ .*

As we showed earlier, the simplest marginals (the zeroth order normalization constant), correspond to the Fourier transform at  $\tau_{(n)}$ , while the first-order marginals correspond to  $\tau_{(n-1,1)}$ , and the second-order unordered marginals correspond to  $\tau_{(n-2,2)}$ . The list goes on and on, with the marginals getting more complicated. At the other end of the spectrum, we have the Fourier coefficients at the representation  $\tau_{(1,1,\dots,1)}$  which exactly recover the original probabilities  $P(\sigma)$ .

We use the word ‘spectrum’ suggestively here, because the different levels of complexity for the marginals are highly reminiscent of the different frequencies for real-valued signals, and a natural question to ask is how the partitions might be ordered with respect to the ‘complexity’ of the corresponding basis functions. In particular how might one characterize this vague notion of complexity for a given partition?

The ‘correct’ characterization, as it turns out, is to use the *dominance ordering* of partitions, which, unlike the ordering on frequencies, is not a linear order, but rather, a partial order.

**Definition 5 (Dominance Ordering)** *Let  $\lambda, \mu$  be partitions of  $n$ . Then  $\lambda \succeq \mu$  (we say  $\lambda$  dominates  $\mu$ ), if for each  $i$ ,  $\sum_{k=1}^i \lambda_k \geq \sum_{k=1}^i \mu_k$ .*

For example,  $(4, 2) \succeq (3, 2, 1)$  since  $4 \geq 3$ ,  $4 + 2 \geq 3 + 2$ , and  $4 + 2 + 0 \geq 3 + 2 + 1$ . However,  $(3, 3)$  and  $(4, 1, 1)$  cannot be compared with respect to the dominance ordering since  $3 \leq 4$ , but  $3 + 3 \geq 4 + 1$ . The ordering over the partitions of  $n = 6$  is depicted in Figure 3(a).

Partitions with fat Ferrers diagrams tend to be greater (with respect to dominance ordering) than those with skinny Ferrers diagrams. Intuitively, representations corresponding to partitions which are *high* in the dominance ordering are ‘low frequency’, while representations corresponding to partitions which are *low* in the dominance ordering are ‘high frequency.’<sup>6</sup>

Having defined a family of intuitive permutation representations over the Symmetric group, we can now ask whether the permutation representations are irreducible or not: the answer in general, is to the negative, due to the fact that it is often possible to reconstruct lower order marginals by summing over the appropriate higher order marginal probabilities. However, it is possible to show

6. The direction of the ordering is slightly counterintuitive given the frequency interpretation, but is standard in the literature.

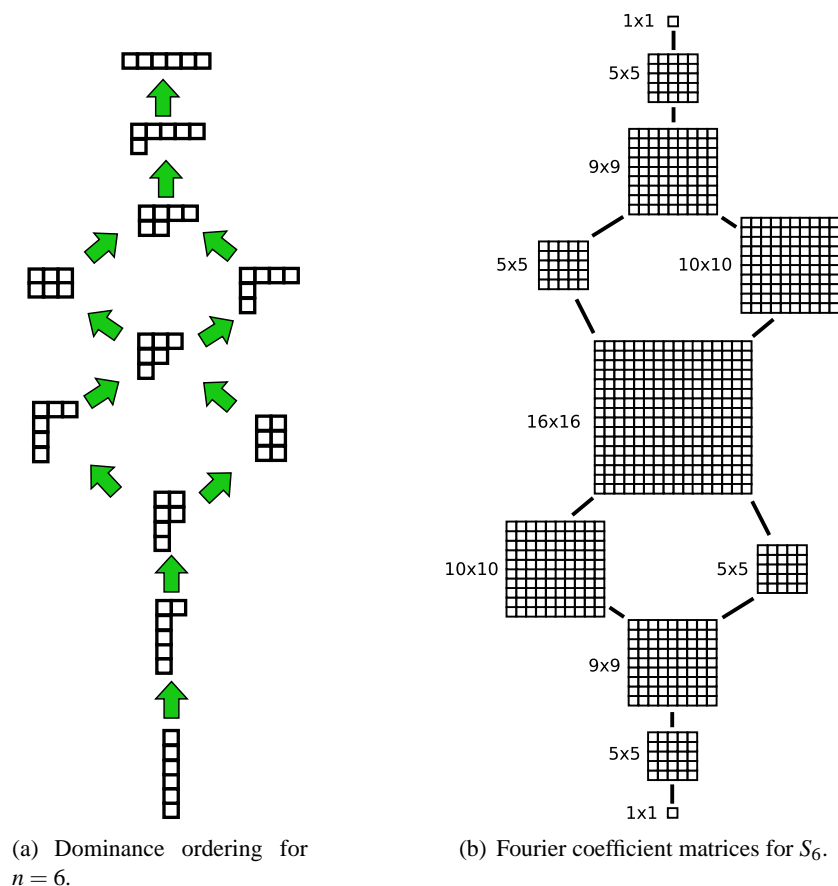


Figure 3: The dominance order for partitions of  $n = 6$  are shown in the left diagram (a). Fat Ferrer’s diagrams tend to be higher in the order and long, skinny diagrams tend to be lower. The corresponding Fourier coefficient matrices for each partition (at irreducible representations) are shown in the right diagram (b). Note that since the Fourier basis functions form a complete basis for the space of functions on the Symmetric group, there must be exactly  $n!$  coefficients in total.

that, for each permutation representation  $\tau_\lambda$ , there exists a corresponding irreducible representation  $\rho_\lambda$ , which, loosely, captures all of the information at the ‘frequency’  $\lambda$  which was not already captured at lower frequency irreducibles. Moreover, it can be shown that there exists no irreducible representation besides those indexed by the partitions of  $n$ . These remarkable results are formalized in the *James Submodule Theorem*, which we state here without proof (see Diaconis 1988, James and Kerber 1981 and Sagan 2001).

**Theorem 6 (James’ Submodule Theorem)**

1. (Uniqueness) For each partition,  $\lambda$ , of  $n$ , there exists an irreducible representation,  $\rho_\lambda$ , which is unique up to equivalence.
2. (Completeness) Every irreducible representation of  $S_n$  corresponds to some partition of  $n$ .

3. There exists a matrix  $C_\lambda$  associated with each partition  $\lambda$ , for which

$$C_\lambda^T \cdot \tau_\lambda(\sigma) \cdot C_\lambda = \bigoplus_{\mu \triangleright \lambda} \bigoplus_{\ell=1}^{K_{\lambda\mu}} \rho_\mu(\sigma), \quad \text{for all } \sigma \in S_n. \quad (5)$$

4.  $K_{\lambda\lambda} = 1$  for all partitions  $\lambda$ .

In plain English, part (3) of the James Submodule theorem says that we can always reconstruct marginal probabilities of  $\lambda$ -tabloids using the Fourier coefficients at irreducibles which lie *at  $\lambda$  and above* in the dominance ordering, if we have knowledge of the matrix  $C_\lambda$  (which can be precomputed using methods detailed in Appendix D), and the multiplicities  $K_{\lambda\mu}$ . In particular, combining Equation 5 with the definition of the Fourier transform, we have that

$$\hat{f}_{\tau_\lambda} = C_\lambda \cdot \left[ \bigoplus_{\mu \triangleright \lambda} \bigoplus_{\ell=1}^{K_{\lambda\mu}} \hat{f}_{\rho_\mu} \right] \cdot C_\lambda^T, \quad (6)$$

and so to obtain marginal probabilities of  $\lambda$ -tabloids, we simply construct a block diagonal matrix using the appropriate irreducible Fourier coefficients, and conjugate by  $C_\lambda$ . The multiplicities  $K_{\lambda\mu}$  are known as the *Kostka numbers* and can be computed using Young's rule (Sagan, 2001). To illustrate using a few examples, we have the following decompositions:

$$\begin{aligned} \tau_{(n)} &\equiv \rho_{(n)}, \\ \tau_{(n-1,1)} &\equiv \rho_{(n)} \oplus \rho_{(n-1,1)}, \\ \tau_{(n-2,2)} &\equiv \rho_{(n)} \oplus \rho_{(n-1,1)} \oplus \rho_{(n-2,2)}, \\ \tau_{(n-2,1,1)} &\equiv \rho_{(n)} \oplus \rho_{(n-1,1)} \oplus \rho_{(n-1,1)} \oplus \rho_{(n-2,2)} \oplus \rho_{(n-2,1,1)}, \\ \tau_{(n-3,3)} &\equiv \rho_{(n)} \oplus \rho_{(n-1,1)} \oplus \rho_{(n-2,2)} \oplus \rho_{(n-3,3)}, \\ \tau_{(n-3,2,1)} &\equiv \rho_{(n)} \oplus \rho_{(n-1,1)} \oplus \rho_{(n-1,1)} \oplus \rho_{(n-2,2)} \oplus \rho_{(n-2,2)} \oplus \rho_{(n-2,1,1)} \oplus \rho_{(n-3,3)} \oplus \rho_{(n-3,2,1)}. \end{aligned}$$

Intuitively, the irreducibles at a partition  $\lambda$  reflect the “pure”  $\lambda^{\text{th}}$ -order effects of the underlying distribution. In other words, the irreducibles at  $\lambda$  form a basis for functions that have “interesting”  $\lambda^{\text{th}}$ -order marginal probabilities, but uniform marginals at all partitions  $\mu$  such that  $\mu \triangleright \lambda$ .

**Example 7** As an example, we demonstrate a “preference” function which is “purely” second-order (unordered) in the sense that its Fourier coefficients are equal to zero at all irreducible representations except  $\rho_{(n-2,2)}$  (and the trivial representation). Consider the function  $f : S_n \rightarrow \mathbb{R}$  defined by:

$$f(\sigma) = \begin{cases} 1 & \text{if } |\sigma(1) - \sigma(2)| \equiv 1 \pmod{n} \\ 0 & \text{otherwise} \end{cases}.$$

Intuitively, imagine seating  $n$  people at a round table with  $n$  chairs, but with the constraint that the first two people, Alice and Bob, are only happy if they are allowed to sit next to each other. In this case,  $f$  can be thought of as the indicator function for the subset of seating arrangements (permutations) which make Alice and Bob happy.

Since  $f$  depends only on the destination of the unordered pair  $\{1, 2\}$ , its Fourier transform is zero at all partitions  $\mu$  such that  $\mu \triangleleft (n-2, 2)$  ( $\hat{f}_\mu = 0$ ). On the other hand, Alice and Bob

$\lambda$	$(n)$	$(n-1, 1)$	$(n-2, 2)$	$(n-2, 1, 1)$	$(n-3, 3)$	$(n-3, 2, 1)$
$\dim \rho_\lambda$	1	$n-1$	$\frac{n(n-3)}{2}$	$\frac{(n-1)(n-2)}{2}$	$\frac{n(n-1)(n-5)}{6}$	$\frac{n(n-2)(n-4)}{3}$

Table 3: Dimensions of low-order irreducible representation matrices.

have no individual preferences for seating, so the first-order “marginals” of  $f$  are uniform, and hence,  $\hat{f}_{(n-1,1)} = 0$ . The Fourier coefficients at irreducibles can be obtained from the second-order (unordered) “marginals” using Equation 5.

$$C_{(n-2,2)}^T \cdot \hat{P}_{\tau_{(n-2,2)}} \cdot C_{(n-2,2)} = \begin{bmatrix} \boxed{Z} & & \\ & \boxed{0} & \\ & & \boxed{\hat{f}_{\rho_{(n-2,2)}}} \end{bmatrix}.$$

The sizes of the irreducible representation matrices are typically much smaller than their corresponding permutation representation matrices. In the case of  $\lambda = (1, \dots, 1)$  for example,  $\dim \tau_\lambda = n!$  while  $\dim \rho_\lambda = 1$ . There is a simple combinatorial algorithm, known as the *Hook Formula* (Sagan, 2001), for computing the dimension of  $\rho_\lambda$ . While we do not discuss it, we provide a few dimensionality computations here (Table 3) to facilitate a discussion of complexity later. Despite providing polynomial sized function approximations, the Fourier coefficient matrices can grow quite fast, and roughly, one would need  $O(n^{2k})$  storage to maintain  $k$ th order marginals. For example, we would need to store  $O(n^8)$  elements to maintain fourth-order marginals. It is worth noting that since the Fourier transform is invertible, there must be  $n!$  Fourier coefficients in total, and so  $\sum_\rho d_\rho^2 = |G| = n!$ . See Figure 3(b) for an example of what the matrices of a complete Fourier transform on  $S_6$  would look like.

In practice, since the irreducible representation matrices are determined only up to equivalence, it is necessary to choose a basis for the irreducible representations in order to explicitly construct the representation matrices. As in Kondor et al. (2007), we use the *Gel’fand-Tsetlin basis* which has several attractive properties, two advantages being that the matrices are real-valued and orthogonal. See Appendix B for details on constructing irreducible matrix representations with respect to the Gel’fand-Tsetlin basis.

### 6. Inference in the Fourier Domain

What we have shown thus far is that there is a principled method for compactly summarizing distributions over permutations based on the idea of bandlimiting—saving only the low-frequency terms of the Fourier transform of a function, which, as we discussed, is equivalent to maintaining a set of low-order marginal probabilities. We now turn to the problem of performing probabilistic inference

using our compact summaries. One of the main advantages of viewing marginals as Fourier coefficients is that it provides a natural principle for formulating polynomial time approximate inference algorithms, which is to rewrite all inference related operations with respect to the Fourier domain, then to perform the Fourier domain operations ignoring high-order terms.

The idea of bandlimiting a distribution is ultimately moot, however, if it becomes necessary to transform back to the primal domain each time an inference operation is called. Naively, the Fourier Transform on  $S_n$  scales as  $O((n!)^2)$ , and even the fastest Fast Fourier Transforms for functions on  $S_n$  are no faster than  $O(n^2 \cdot n!)$  (see Maslen 1998 for example). To resolve this issue, we present a formulation of inference which operates solely in the Fourier domain, allowing us to avoid a costly transform. We begin by discussing exact inference in the Fourier domain, which is no more tractable than the original problem because there are  $n!$  Fourier coefficients, but it will allow us to discuss the bandlimiting approximation in the next section. There are two operations to consider: prediction/rollup, and conditioning. While we have motivated both of these operations in the familiar context of hidden Markov models, they are fundamental and appear in many other settings. The assumption for the rest of this section is that the Fourier transforms of the transition and observation models are known. We discuss methods for obtaining the models in Section 8. The main results of this section (excluding the discussions about complexity) extend naturally to other finite groups besides  $S_n$ .

### 6.1 Fourier Prediction/Rollup

We will consider one particular class of transition models—that of random walks over a group, which assumes that  $\sigma^{(t+1)}$  is generated from  $\sigma^{(t)}$  by drawing a random permutation  $\pi^{(t)}$  from some distribution  $Q^{(t)}$  and setting  $\sigma^{(t+1)} = \pi^{(t)}\sigma^{(t)}$ .<sup>7</sup> In our identity management example,  $\pi^{(t)}$  represents a random identity permutation that might occur among tracks when they get close to each other (what we call a *mixing event*). For example,  $Q(1,2) = 1/2$  means that Tracks 1 and 2 swapped identities with probability 1/2. The random walk model also appears in many other applications such as modeling card shuffles (Diaconis, 1988).

The motivation behind the random walk transition model is that it allows us to write the prediction/rollup operation as a *convolution* of distributions on a group. The extension of the familiar notion of convolution to groups simply replaces additions and subtractions by analogous group operations (function composition and inverse, respectively):

**Definition 7** Let  $Q$  and  $P$  be probability distributions on a group  $G$ . Define the convolution<sup>8</sup> of  $Q$  and  $P$  to be the function  $[Q * P](\sigma_1) = \sum_{\sigma_2} Q(\sigma_1\sigma_2^{-1})P(\sigma_2)$ .

Using Definition 7, we see that the prediction/rollup step can be written as:

$$\begin{aligned} P(\sigma^{(t+1)}) &= \sum_{\sigma^{(t)}} P(\sigma^{(t+1)} | \sigma^{(t)}) \cdot P(\sigma^{(t)}), \\ &= \sum_{\{\sigma^{(t)}, \pi^{(t)} : \sigma^{(t+1)} = \pi^{(t)} \cdot \sigma^{(t)}\}} Q^{(t)}(\pi^{(t)}) \cdot P(\sigma^{(t)}), \end{aligned}$$

7. We place  $\pi$  on the left side of the multiplication because we want it to permute tracks and not identities. Had we defined  $\pi$  to map from tracks to identities (instead of identities to tracks), then  $\pi$  would be multiplied from the right. Besides left versus right multiplication, there are no differences between the two conventions.

8. Note that this definition of convolution on groups is *strictly* a generalization of convolution of functions on the real line, and is a non-commutative operation for non-Abelian groups. Thus the distribution  $P * Q$  is not necessarily the same as  $Q * P$ .

$$\begin{aligned}
 & \text{(Right-multiplying both sides of } \sigma^{(t+1)} = \pi^{(t)} \sigma^{(t)} \\
 & \text{by } (\sigma^{(t)})^{-1}, \text{ we see that } \pi^{(t)} \text{ can be replaced by } \sigma^{(t+1)}(\sigma^{(t)})^{-1}), \\
 & = \sum_{\sigma^{(t)}} Q^{(t)}(\sigma^{(t+1)} \cdot (\sigma^{(t)})^{-1}) \cdot P(\sigma^{(t)}), \\
 & = [Q^{(t)} * P](\sigma^{(t+1)}).
 \end{aligned}$$

As with Fourier transforms on the real line, the Fourier coefficients of the convolution of distributions  $P$  and  $Q$  on groups can be obtained from the Fourier coefficients of  $P$  and  $Q$  individually, using the *convolution theorem* (see also Diaconis 1988):

**Proposition 8 (Convolution Theorem)** *Let  $Q$  and  $P$  be probability distributions on a group  $G$ . For any representation  $\rho$ ,*

$$[\widehat{Q * P}]_{\rho} = \widehat{Q}_{\rho} \cdot \widehat{P}_{\rho},$$

where the operation on the right side is matrix multiplication.

Therefore, assuming that the Fourier transforms  $\widehat{P}_{\rho}^{(t)}$  and  $\widehat{Q}_{\rho}^{(t)}$  are given, the prediction/rollup update rule is simply:

$$\widehat{P}_{\rho}^{(t+1)} \leftarrow \widehat{Q}_{\rho}^{(t)} \cdot \widehat{P}_{\rho}^{(t)}.$$

Note that the update only requires knowledge of  $\widehat{P}$  and does not require  $P$ . Furthermore, the update is *pointwise* in the Fourier domain in the sense that the coefficients at the representation  $\rho$  affect  $\widehat{P}_{\rho}^{(t+1)}$  only at  $\rho$ . Consequently, prediction/rollup updates in the Fourier domain never increase the representational complexity. For example, if we maintain third-order marginals, then a single step of prediction/rollup called at time  $t$  returns the *exact* third-order marginals at time  $t + 1$ , and nothing more.

**Example 8** *We run the prediction/rollup routines on the first two time steps of the example in Figure 2, first in the primal domain, then in the Fourier domain. At each mixing event, two tracks,  $i$  and  $j$ , swap identities with some probability. Using a mixing model given by:*

$$Q(\pi) = \begin{cases} 3/4 & \text{if } \pi = \varepsilon \\ 1/4 & \text{if } \pi = (i, j) \\ 0 & \text{otherwise} \end{cases},$$

we obtain results shown in Tables 4 and 5.

### 6.1.1 COMPLEXITY OF PREDICTION/ROLLUP

We will discuss complexity in terms of the dimension of the largest maintained irreducible Fourier coefficient matrix, which we will denote by  $d_{max}$  (see Table 3 for irreducible dimensions). If we maintain  $2^{nd}$  order marginals, for example, then  $d_{max} = O(n^2)$ , and if we maintain  $3^{rd}$  order marginals, then  $d_{max} = O(n^3)$ .

Performing a single prediction/rollup step in the Fourier domain involves performing a single matrix multiplication for each irreducible and thus requires  $O(d_{max}^3)$  time using the naive multiplication algorithm.

$\sigma$	$P^{(0)}$	$Q^{(1)}$	$P^{(1)}$	$Q^{(2)}$	$P^{(2)}$
$\varepsilon$	1	3/4	3/4	3/4	9/16
(1,2)	0	1/4	1/4	0	3/16
(2,3)	0	0	0	0	0
(1,3)	0	0	0	1/4	3/16
(1,2,3)	0	0	0	0	1/16
(1,3,2)	0	0	0	0	0

Table 4: Primal domain prediction/rollup example.

	$\hat{P}^{(0)}$	$\hat{Q}^{(1)}$	$\hat{P}^{(1)}$	$\hat{Q}^{(2)}$	$\hat{P}^{(2)}$
$\rho_{(3)}$	1	1	1	1	1
$\rho_{(2,1)}$	$\begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$	$\begin{bmatrix} 1/2 & 0 \\ 0 & 1 \end{bmatrix}$	$\begin{bmatrix} 1/2 & 0 \\ 0 & 1 \end{bmatrix}$	$\begin{bmatrix} 7/8 & -\sqrt{3}/8 \\ -\sqrt{3}/8 & 5/8 \end{bmatrix}$	$\begin{bmatrix} 7/16 & -\sqrt{3}/8 \\ -\sqrt{3}/16 & 5/8 \end{bmatrix}$
$\rho_{(1,1,1)}$	1	1/2	1/2	1/2	1/4

Table 5: Fourier domain prediction/rollup example.

In certain situations, faster updates can be achieved. For example, in the pairwise mixing model of Example 8, the Fourier transform of  $Q$  distribution takes the form:  $\hat{Q}_{\rho_\lambda} = \alpha I_{d_\lambda} + \beta \rho_\lambda(i, j)$ , where  $I_{d_\lambda}$  is the  $d_\lambda \times d_\lambda$  identity matrix (see also Section 8). As it turns out, the matrix  $\rho_\lambda(i, j)$  can be factored into a product of  $O(n)$  sparse matrices each with at most  $O(d_\lambda)$  nonzero entries. To see why, recall the elementary fact that the transposition  $(i, j)$  factors into a sequence of  $O(n)$  adjacent transpositions:

$$(i, j) = (i, i+1)(i+1, i+2) \cdots (j-1, j)(j-2, j-1) \cdots (i+1, i+2)(i, i+1).$$

If we use the Gel'fand-Tsetlin basis adapted to the subgroup chain  $S_1 \subset \cdots \subset S_n$  (see Appendix B), then we also know that the irreducible representation matrices evaluated at adjacent transpositions are sparse with no more than  $O(d_{\max}^2)$  nonzero entries. Thus by carefully exploiting sparsity during the prediction/rollup algorithm, one can achieve an  $O(nd_{\max}^2)$  update, which is faster than  $O(d_{\max}^3)$  as long as one uses more than first-order terms.

### 6.1.2 LIMITATIONS OF RANDOM WALK MODELS

While the random walk assumption captures a rather general family of transition models, there do exist certain models which cannot be written as a random walk on a group. In particular, one limitation is that the prediction/rollup update for a random walk model can only increase the entropy of the distribution. As with Kalman filters, localization is thus impossible without making obser-

vations.<sup>9</sup> Shin et al. (2005) show that the entropy must increase for a certain kind of random walk on  $S_n$  (where  $\pi$  could be either the identity or the transposition  $(i, j)$ ), but in fact, the result is easily generalized for any random walk mixing model and for any finite group.

**Proposition 9**

$$H [P^{(t+1)}(\sigma^{(t+1)})] \geq \max \left\{ H [Q^{(t)}(\tau^{(t)})], H [P^{(t)}(\sigma^{(t)})] \right\},$$

where  $H [P(\sigma)]$  denotes the statistical entropy functional,  $H[P(\sigma)] = -\sum_{\sigma \in G} P(\sigma) \log P(\sigma)$ .

**Proof** We have:

$$\begin{aligned} P^{(t+1)}(\sigma^{(t+1)}) &= [Q^{(t)} * P^{(t)}] (\sigma^{(t+1)}) \\ &= \sum_{\sigma^{(t)}} Q(\sigma^{(t+1)} \cdot (\sigma^{(t)})^{-1}) P^{(t)}(\sigma^{(t)}) \end{aligned}$$

Applying the Jensen Inequality to the entropy function (which is concave) yields:

$$\begin{aligned} H [P^{(t+1)}(\sigma^{(t+1)})] &\geq \sum_{\sigma^{(t)}} P^{(t)}(\sigma^{(t)}) H [Q^{(t)}(\sigma \cdot (\sigma^{(t)})^{-1})], && \text{(Jensen's inequality)} \\ &= \sum_{\sigma^{(t)}} P^{(t)}(\sigma^{(t)}) H [Q^{(t)}(\sigma)], && \text{(translation invariance of entropy)} \\ &= H [Q^{(t)}(\sigma)], && \text{(since } \sum_{\sigma^{(t)}} P^{(t)}(\sigma^{(t)}) = 1). \end{aligned}$$

The proof that  $H [P^{(t+1)}(\sigma^{(t+1)})] \geq H [P^{(t)}(\sigma^{(t)})]$  is similar with the exception that we must rewrite the convolution so that the sum ranges over  $\tau^{(t)}$ .

$$\begin{aligned} P^{(t+1)}(\sigma^{(t+1)}) &= [Q^{(t)} * P^{(t)}] (\sigma^{(t+1)}), \\ &= \sum_{\tau^{(t)}} Q^{(t)}(\tau^{(t)}) P^{(t)}((\tau^{(t)})^{-1} \cdot \sigma^{(t+1)}). \end{aligned}$$

■

**Example 9** This example is based on one from Diaconis (1988). Consider a deck of cards numbered  $\{1, \dots, n\}$ . Choose a random permutation of cards by first picking two cards independently, and swapping (a card might be swapped with itself), yielding the following probability distribution over  $S_n$ :

$$Q(\pi) = \begin{cases} \frac{1}{n} & \text{if } \pi = \epsilon \\ \frac{2}{n^2} & \text{if } \pi \text{ is a transposition} \\ 0 & \text{otherwise} \end{cases} \quad (7)$$

---

9. In general, if we are not constrained to using linear Gaussian models, it is possible to localize with no observations. Consider a robot walking along the unit interval on the real line (which is not a group). If the position of the robot is unknown, one easy localization strategy might be to simply drive the robot to the right, with the knowledge that given ample time, the robot will slam into the ‘wall’, at which point it will have been localized. With random walk based models on groups however, these strategies are impossible—imagine the same robot walking around the unit circle—since, in some sense, the group structure prevents the existence of ‘walls’.



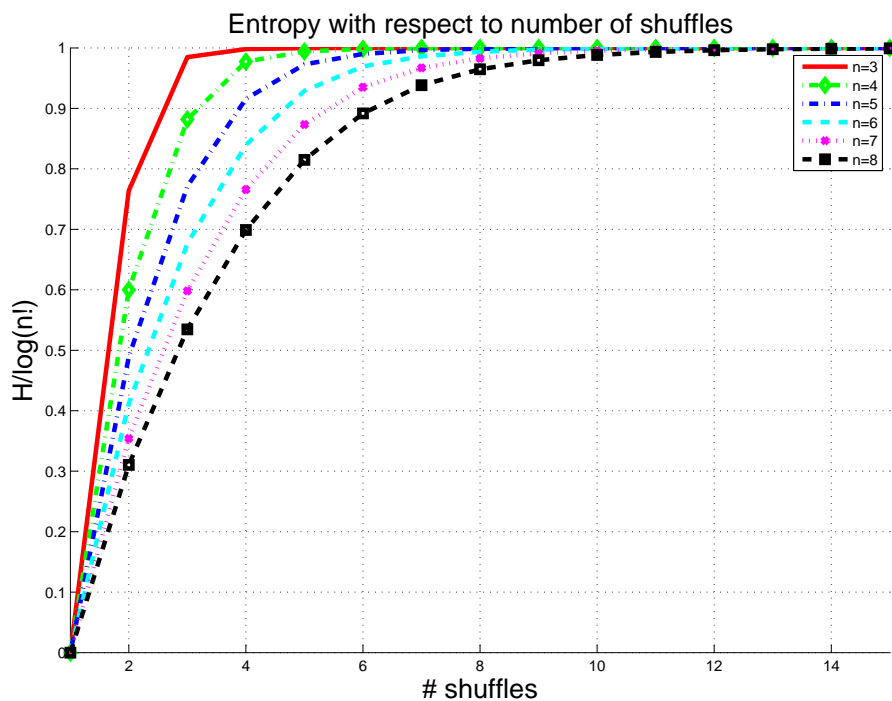


Figure 4: We start with a deck of cards in sorted order, and perform fifteen consecutive shuffles according to the rule given in Equation 7. The plot shows the entropy of the distribution over permutations with respect to the number of shuffles for  $n = 3, 4, \dots, 8$ . When  $H(P)/\log(n!) = 1$ , the distribution has become uniform.

*Repeating the above process for generating random permutations  $\pi$  gives a transition model for a hidden Markov model over the symmetric group. We can also see (Figure 4) that the entropy of the deck increases monotonically with each shuffle, and that repeated shuffles with  $Q(\pi)$  eventually bring the deck to the uniform distribution.*

## 6.2 Fourier Conditioning

In contrast with the prediction/rollup operation, conditioning can potentially increase the representational complexity. As an example, suppose that we know the following first-order marginal probabilities:

$$P(\text{Alice is at Track 1 or Track 2}) = .9, \text{ and}$$

$$P(\text{Bob is at Track 1 or Track 2}) = .9.$$

If we then make the following first-order observation:

$$P(\text{Cathy is at Track 1 or Track 2}) = 1,$$

then it can be inferred that Alice and Bob cannot *both* occupy Tracks 1 and 2 at the same time, that is,

$$P(\{\text{Alice,Bob}\} \text{ occupy Tracks } \{1,2\}) = 0,$$

demonstrating that after conditioning, we are left with knowledge of second-order (unordered) marginals despite the fact that the prior and likelihood functions were only known up to first-order. Intuitively, the example shows that conditioning “smears” information from low-order Fourier coefficients to high-order coefficients, and that one cannot hope for a pointwise operation as was afforded by prediction/rollup. We now show precisely how irreducibles of different complexities “interact” with each other in the Fourier domain during conditioning.

An application of Bayes rule to find a posterior distribution  $P(\sigma|z)$  after observing some evidence  $z$  requires two steps: a *pointwise product* of likelihood  $P(z|\sigma)$  and prior  $P(\sigma)$ , followed by a normalization step:

$$P(\sigma|z) = \eta \cdot P(z|\sigma) \cdot P(\sigma).$$

For notational convenience, we will refer to the likelihood function as  $L(z|\sigma)$  henceforth. We showed earlier that the normalization constant  $\eta^{-1} = \sum_{\sigma} L(z|\sigma) \cdot P(\sigma)$  is given by the Fourier transform of  $\widehat{L^{(t)}P^{(t)}}$  at the trivial representation—and therefore the normalization step of conditioning can be implemented by simply dividing each Fourier coefficient by the scalar  $\left[\widehat{L^{(t)}P^{(t)}}\right]_{\rho_{(n)}}$ .

The pointwise product of two functions  $f$  and  $g$ , however, is trickier to formulate in the Fourier domain. For functions on the real line, the pointwise product of functions can be implemented by convolving the Fourier coefficients of  $\hat{f}$  and  $\hat{g}$ , and so a natural question is: can we apply a similar operation for functions over general groups? Our answer to this is that there is an analogous (but more complicated) notion of convolution in the Fourier domain of a general finite group. We present a convolution-based conditioning algorithm which we call *Kronecker Conditioning*, which, in contrast to the pointwise nature of the Fourier Domain prediction/rollup step, and much like convolution, smears the information at an irreducible  $\rho_v$  to other irreducibles.

### 6.2.1 FOURIER TRANSFORM OF THE POINTWISE PRODUCT

Our approach to computing the Fourier transform of the pointwise product in terms of  $\hat{f}$  and  $\hat{g}$  is to manipulate the function  $f(\sigma)g(\sigma)$  so that it can be seen as the result of an inverse Fourier transform (Equation 3). Hence, the goal will be to find matrices  $R_v$  (as a function of  $\hat{f}, \hat{g}$ ) such that for any  $\sigma \in G$ ,

$$f(\sigma) \cdot g(\sigma) = \frac{1}{|G|} \sum_v d_{\rho_v} \text{Tr}(R_v^T \cdot \rho_v(\sigma)), \tag{8}$$

after which we will be able to read off the Fourier transform of the pointwise product as  $\left[\widehat{fg}\right]_{\rho_v} = R_v$ .

For any  $\sigma \in G$ , we can write the pointwise product in terms of  $\hat{f}$  and  $\hat{g}$  using the inverse Fourier transform:

$$\begin{aligned} f(\sigma) \cdot g(\sigma) &= \left[ \frac{1}{|G|} \sum_{\lambda} d_{\rho_{\lambda}} \text{Tr}(\hat{f}_{\rho_{\lambda}}^T \cdot \rho_{\lambda}(\sigma)) \right] \cdot \left[ \frac{1}{|G|} \sum_{\mu} d_{\rho_{\mu}} \text{Tr}(\hat{g}_{\rho_{\mu}}^T \cdot \rho_{\mu}(\sigma)) \right] \\ &= \left( \frac{1}{|G|} \right)^2 \sum_{\lambda, \mu} d_{\rho_{\lambda}} d_{\rho_{\mu}} \left[ \text{Tr}(\hat{f}_{\rho_{\lambda}}^T \cdot \rho_{\lambda}(\sigma)) \cdot \text{Tr}(\hat{g}_{\rho_{\mu}}^T \cdot \rho_{\mu}(\sigma)) \right]. \end{aligned} \tag{9}$$

<ol style="list-style-type: none"> <li>1. If <math>A</math> and <math>B</math> are square, <math>\text{Tr}(A \otimes B) = (\text{Tr}A) \cdot (\text{Tr}B)</math>.</li> <li>2. <math>(A \otimes B) \cdot (C \otimes D) = AC \otimes BD</math>.</li> <li>3. Let <math>A</math> be an <math>n \times n</math> matrix, and <math>C</math> an invertible <math>n \times n</math> matrix. Then <math>\text{Tr}A = \text{Tr}(C^{-1}AC)</math>.</li> <li>4. Let <math>A</math> be an <math>n \times n</math> matrix and <math>B_i</math> be matrices of size <math>m_i \times m_i</math> where <math>\sum_i m_i = n</math>. Then <math>\text{Tr}(A \cdot (\bigoplus_i B_i)) = \sum_i \text{Tr}(A_i \cdot B_i)</math>, where <math>A_i</math> is the block of <math>A</math> corresponding to block <math>B_i</math> in the matrix <math>(\bigoplus_i B_i)</math>.</li> </ol>
--

Table 6: Matrix Identities used in Proposition 10.

Now we want to manipulate this product of traces in the last line to be just one trace (as in Equation 8), by appealing to some properties of the *Kronecker Product*. The Kronecker product of an  $n \times n$  matrix  $U = (u_{i,j})$  by an  $m \times m$  matrix  $V$ , is defined to be the  $nm \times nm$  matrix

$$U \otimes V = \begin{pmatrix} u_{1,1}V & u_{1,2}V & \dots & u_{1,n}V \\ u_{2,1}V & u_{2,2}V & \dots & u_{2,n}V \\ \vdots & \vdots & \ddots & \vdots \\ u_{n,1}V & u_{n,2}V & \dots & u_{n,n}V \end{pmatrix}.$$

We summarize some important matrix properties in Table 6. The connection to our problem is given by matrix property 1. Applying this to Equation 9, we have:

$$\begin{aligned} \text{Tr}\left(\hat{f}_{\rho_\lambda}^T \cdot \rho_\lambda(\sigma)\right) \cdot \text{Tr}\left(\hat{g}_{\rho_\mu}^T \cdot \rho_\mu(\sigma)\right) &= \text{Tr}\left(\left(\hat{f}_{\rho_\lambda}^T \cdot \rho_\lambda(\sigma)\right) \otimes \left(\hat{g}_{\rho_\mu}^T \cdot \rho_\mu(\sigma)\right)\right) \\ &= \text{Tr}\left(\left(\hat{f}_{\rho_\lambda} \otimes \hat{g}_{\rho_\mu}\right)^T \cdot (\rho_\lambda(\sigma) \otimes \rho_\mu(\sigma))\right), \end{aligned}$$

where the last line follows by Property 2. The term on the left,  $\hat{f}_{\rho_\lambda} \otimes \hat{g}_{\rho_\mu}$ , is a matrix of coefficients. The term on the right,  $\rho_\lambda(\sigma) \otimes \rho_\mu(\sigma)$ , itself happens to be a representation, called the *Kronecker (or Tensor) Product Representation*. In general, the Kronecker product representation is reducible, and so it can be decomposed into a direct sum of irreducibles. In particular, if  $\rho_\lambda$  and  $\rho_\mu$  are any two irreducibles of  $G$ , there exists a similarity transform  $C_{\lambda\mu}$  such that, for any  $\sigma \in G$ ,

$$C_{\lambda\mu}^{-1} \cdot [\rho_\lambda \otimes \rho_\mu](\sigma) \cdot C_{\lambda\mu} = \bigoplus_{\mathbf{v}} \bigoplus_{\ell=1}^{z_{\lambda\mu\mathbf{v}}} \rho_{\mathbf{v}}(\sigma). \quad (10)$$

The  $\bigoplus$  symbols here refer to a matrix direct sum as in Equation 2,  $\mathbf{v}$  indexes over all irreducible representations of  $S_n$ , while  $\ell$  indexes over a number of *copies* of  $\rho_{\mathbf{v}}$  which appear in the decomposition. We index blocks on the right side of this equation by pairs of indices  $(\mathbf{v}, \ell)$ . The number of copies of each  $\rho_{\mathbf{v}}$  (for the tensor product pair  $\rho_\lambda \otimes \rho_\mu$ ) is denoted by the integer  $z_{\lambda\mu\mathbf{v}}$ , the collection of which, taken over all triples  $(\lambda, \mu, \mathbf{v})$ , are commonly referred to as the *Clebsch-Gordan series*. Note that we allow the  $z_{\lambda\mu\mathbf{v}}$  to be zero, in which case  $\rho_{\mathbf{v}}$  does not contribute to the direct sum. The matrices  $C_{\lambda\mu}$  are known as the *Clebsch-Gordan coefficients*. The *Kronecker Product Decomposition* problem is that of finding the irreducible components of the Kronecker product representation, and thus to find the Clebsch-Gordan series/coefficients for each pair of irreducible representations  $(\rho_\lambda, \rho_\mu)$ .

Decomposing the Kronecker product inside Equation 10 using the Clebsch-Gordan series and coefficients yields the desired Fourier transform, which we summarize in the form of a proposition. In the case that  $f$  and  $g$  are defined over an Abelian group, we will show that the following formulas reduce to the familiar form of convolution.

**Proposition 10** *Let  $\hat{f}, \hat{g}$  be the Fourier transforms of functions  $f$  and  $g$  respectively, and for each ordered pair of irreducibles  $(\rho_\lambda, \rho_\mu)$ , define:  $A_{\lambda\mu} \triangleq C_{\lambda\mu}^{-1} \cdot (\hat{f}_{\rho_\lambda} \otimes \hat{g}_{\rho_\mu}) \cdot C_{\lambda\mu}$ . Then the Fourier transform of the pointwise product  $fg$  is:*

$$[\widehat{fg}]_{\rho_\nu} = \frac{1}{d_{\rho_\nu}|G|} \sum_{\lambda\mu} d_{\rho_\lambda} d_{\rho_\mu} \sum_{\ell=1}^{z_{\lambda\mu\nu}} A_{\lambda\mu}^{(\nu, \ell)}, \quad (11)$$

where  $A_{\lambda\mu}^{(\nu, \ell)}$  is the block of  $A_{\lambda\mu}$  corresponding to the  $(\nu, \ell)$  block in  $\bigoplus_\nu \bigoplus_{\ell=1}^{z_{\lambda\mu\nu}} \rho_\nu$  from Equation 10.

**Proof** We use the fact that  $C_{\lambda\mu}$  is an orthogonal matrix for all pairs  $(\rho_\lambda, \rho_\mu)$ , that is,  $C_{\lambda\mu}^T \cdot C_{\lambda\mu} = I$ .

$$\begin{aligned} f(\sigma) \cdot g(\sigma) &= \left[ \frac{1}{|G|} \sum_{\lambda} d_{\rho_\lambda} \text{Tr} \left( \hat{f}_{\rho_\lambda}^T \cdot \rho_\lambda(\sigma) \right) \right] \cdot \left[ \frac{1}{|G|} \sum_{\mu} d_{\rho_\mu} \text{Tr} \left( \hat{g}_{\rho_\mu}^T \cdot \rho_\mu(\sigma) \right) \right] \\ &= \left( \frac{1}{|G|} \right)^2 \sum_{\lambda, \mu} d_{\rho_\lambda} d_{\rho_\mu} \left[ \text{Tr} \left( \hat{f}_{\rho_\lambda}^T \cdot \rho_\mu(\sigma) \right) \cdot \text{Tr} \left( \hat{g}_{\rho_\mu}^T \cdot \rho_\mu(\sigma) \right) \right] \\ \text{(by Property 1)} &= \left( \frac{1}{|G|} \right)^2 \sum_{\lambda, \mu} d_{\rho_\lambda} d_{\rho_\mu} \left[ \text{Tr} \left( \left( \hat{f}_{\rho_\lambda}^T \cdot \rho_\lambda(\sigma) \right) \otimes \left( \hat{g}_{\rho_\mu}^T \cdot \rho_\mu(\sigma) \right) \right) \right] \\ \text{(by Property 2)} &= \left( \frac{1}{|G|} \right)^2 \sum_{\lambda, \mu} d_{\rho_\lambda} d_{\rho_\mu} \text{Tr} \left( \left( \hat{f}_{\rho_\lambda} \otimes \hat{g}_{\rho_\mu} \right)^T \cdot \left( \rho_\lambda(\sigma) \otimes \rho_\mu(\sigma) \right) \right) \\ \text{(by Property 3)} &= \left( \frac{1}{|G|} \right)^2 \sum_{\lambda, \mu} d_{\rho_\lambda} d_{\rho_\mu} \text{Tr} \left( C_{\lambda\mu}^T \cdot \left( \hat{f}_{\rho_\lambda} \otimes \hat{g}_{\rho_\mu} \right)^T \cdot C_{\lambda\mu} \right. \\ &\quad \left. \cdot C_{\lambda\mu}^T \cdot \left( \rho_\lambda(\sigma) \otimes \rho_\mu(\sigma) \right) \cdot C_{\lambda\mu} \right) \\ \text{(by definition of } C_{\lambda\mu} \text{ and } A_{\lambda\mu}) &= \left( \frac{1}{|G|} \right)^2 \sum_{\lambda, \mu} d_{\rho_\lambda} d_{\rho_\mu} \text{Tr} \left( A_{\lambda\mu}^T \cdot \left( \bigoplus_{\nu} \bigoplus_{\ell=1}^{z_{\lambda\mu\nu}} \rho_\nu(\sigma) \right) \right) \\ \text{(by Property 4)} &= \frac{1}{|G|^2} \sum_{\lambda\mu} d_{\rho_\lambda} d_{\rho_\mu} \sum_{\nu} d_{\rho_\nu} \sum_{\ell=1}^{z_{\lambda\mu\nu}} \text{Tr} \left( \left( d_{\rho_\nu}^{-1} A_{\lambda\mu}^{(\nu, \ell)} \right)^T \rho_\nu(\sigma) \right) \\ \text{(rearranging terms)} &= \frac{1}{|G|} \sum_{\nu} d_{\rho_\nu} \text{Tr} \left[ \left( \sum_{\lambda\mu} \sum_{\ell=1}^{z_{\lambda\mu\nu}} \frac{d_{\rho_\lambda} d_{\rho_\mu}}{d_{\rho_\nu} |G|} A_{\lambda\mu}^{(\nu, \ell)} \right)^T \rho_\nu(\sigma) \right]. \end{aligned}$$

Recognizing the last expression as an inverse Fourier transform completes the proof. ■

The Clebsch-Gordan series,  $z_{\lambda\mu\nu}$ , plays an important role in Equation 11, which says that the  $(\rho_\lambda, \rho_\mu)$  cross-term contributes to the pointwise product at  $\rho_\nu$  *only* when  $z_{\lambda\mu\nu} > 0$ . In the simplest

case, we have that

$$z_{(n),\mu,\nu} = \begin{cases} 1 & \text{if } \mu = \nu \\ 0 & \text{otherwise} \end{cases},$$

which is true since  $\rho_{(n)}(\sigma) = 1$  for all  $\sigma \in S_n$ . As another example, it is known that:

$$\rho_{(n-1,1)} \otimes \rho_{(n-1,1)} \equiv \rho_{(n)} \oplus \rho_{(n-1,1)} \oplus \rho_{(n-2,2)} \oplus \rho_{(n-2,1,1)},$$

or equivalently,

$$z_{(n-1,1),(n-1,1),\nu} = \begin{cases} 1 & \text{if } \nu \text{ is one of } (n), (n-1,1), (n-2,2), \text{ or } (n-2,1,1) \\ 0 & \text{otherwise} \end{cases}.$$

So if the Fourier transforms of the likelihood and prior are zero past the first two irreducibles ( $(n)$  and  $(n-1,1)$ ), then a single conditioning step results in a Fourier transform which, in general, carries second-order information at  $(n-2,2)$  and  $(n-2,1,1)$ , but is guaranteed to be zero past the first four irreducibles  $(n)$ ,  $(n-1,1)$ ,  $(n-2,2)$  and  $(n-2,1,1)$ .

As far as we know, there are no analytical formulas for finding the entire Clebsch-Gordan series or coefficients, and in practice, acquiring the coefficients requires considerable precomputation. We emphasize however, that as fundamental constants related to the irreducibles of the Symmetric group, they need only be computed *once and for all* (like the digits of  $\pi$ , for example) and can be stored in a table for all future reference. For a detailed discussion of techniques for computing the Clebsch-Gordan series/coefficients, see Appendix D. We have made a set of precomputed coefficients available on our lab website,<sup>10</sup> but we will assume throughout the rest of the paper that both the series and coefficients have been made available as a lookup table.

As a final remark, note that Proposition 10 can be rewritten somewhat more intuitively by absorbing the scalars and submatrices of the Clebsch-Gordan coefficients into projection matrices  $P_{\lambda\mu}^{(\nu,\ell)}$ .

**Proposition 11** *Let  $\hat{f}, \hat{g}$  be the Fourier transforms of functions  $f$  and  $g$  respectively. For each triple of partitions  $(\lambda, \mu, \nu)$  there exists a positive integer  $z_{\lambda,\mu,\nu}$  and projection operators  $P_{\lambda\mu}^{(\nu,\ell)}$  for each  $\ell \in \{1, 2, \dots, z_{\lambda,\mu,\nu}\}$  such that the Fourier transform of the pointwise product  $fg$  is:*

$$\left[ \widehat{fg} \right]_{\rho_\nu} = \sum_{\lambda,\mu} \sum_{\ell=1}^{z_{\lambda,\mu,\nu}} (P_{\lambda\mu}^{(\nu,\ell)})^T \cdot (\hat{f}_{\rho_\lambda} \otimes \hat{g}_{\rho_\mu}) \cdot P_{\lambda\mu}^{(\nu,\ell)}. \quad (12)$$

When  $f$  and  $g$  are functions on an Abelian group  $G$ , then it is a well known fact that all irreducible representations are one-dimensional, and so Equation 12 reduces to  $\left[ \widehat{fg} \right]_{\rho_\nu} = \sum_{\lambda,\mu} (\hat{f}_{\rho_\lambda} \cdot \hat{g}_{\rho_\mu})$ , where all the tensor products have simply become scalar multiplications and the familiar definition of convolution is recovered.

### 6.2.2 COMPLEXITY OF CONDITIONING

The complexity of (bandlimited) conditioning (assuming precomputed Clebsch-Gordan series/coefficients) depends on the order of the coefficients maintained for both the prior and the

10. See <http://www.select.cs.cmu.edu/data/index.html>.

observation model. However, it is difficult to state a general complexity bound for arbitrary finite groups due to our limited understanding of the Clebsch-Gordan series. Here we consider conditioning only on the symmetric group of order  $n$  with the assumption that the number of irreducibles maintained is very small (and in particular, not allowed to grow with respect to  $n$ ). Our assumption is realistic in practice since for moderately large  $n$ , it is impractical to consider maintaining higher than, say, third-order terms. If we denote the dimension of the largest maintained irreducibles of the prior and likelihood by  $d_{max}^{prior}$  and  $d_{max}^{obs}$ , respectively, then the complexity of conditioning is dominated by the step that forms a matrix  $C^T \cdot (A \otimes B) \cdot C$ , where the matrices  $A \otimes B$  and  $C$  are each  $(d_{max}^{prior} \cdot d_{max}^{obs})$ -dimensional. Note, however, that since we are only interested in certain blocks of  $C^T \cdot (A \otimes B) \cdot C$ , the full matrix need not be computed. In particular, the largest extracted block has size  $d_{max}^{prior}$ , and so the complexity of conditioning is  $O\left((d_{max}^{obs})^2 (d_{max}^{prior})^3\right)$  using the naive matrix multiplication algorithm.

In some situations (see Section 8), the observation model is fully specified by first-order Fourier terms. In such cases,  $d_{max}^{obs} = O(n)$  and we can perform conditioning in the Fourier domain in  $O(n^2 \cdot (d_{max}^{prior})^3)$  time. If a model is fully specified by second-order terms, for example, then the update requires  $O(n^4 \cdot (d_{max}^{prior})^3)$  time.

To speed up conditioning, one can often exploit matrix sparsity in two ways. First, we observe that the Clebsch-Gordan coefficient matrices are often sparse (we cannot yet prove this, see Figure 10.1) and so we can save a conjectured factor of  $(d_{max}^{prior} \cdot d_{max}^{obs})$  in practice. Secondly, for certain coset-based observation models (see Section 8), we can show that (under an appropriate relabeling of identities and tracks), the Fourier coefficient matrices of the observation model are sparse (with  $O(d_{max}^{obs})$  or sometimes even  $O(1)$  nonzero entries for  $\hat{L}_\lambda$ ). For the simplest observations which take the form (“Identity  $j$  is at track  $j$ ”), for example, we can obtain  $O((d_{max}^{prior})^3)$  running time (without accounting for the conjectured sparsity of the Clebsch-Gordan coefficients), which matches the time required for the prediction/rollup update. See Appendix B for details.

We now conclude our section on inference with a fully worked example of Kronecker conditioning.

**Example 10** *For this example, refer to Table 2 for the representations of  $S_3$ . Given functions  $f, g : S_3 \rightarrow \mathbb{R}$ , we will compute the Fourier transform of the pointwise product  $f \cdot g$ .*

*Since there are three irreducibles, there are nine tensor products  $\rho_\lambda \otimes \rho_\mu$  to decompose, six of which are trivial either because they are one-dimensional, or involve tensoring against the trivial representation. The nontrivial tensor products to consider are  $\rho_{(2,1)} \otimes \rho_{(1,1,1)}$ ,  $\rho_{(1,1,1)} \otimes \rho_{(2,1)}$  and  $\rho_{(2,1)} \otimes \rho_{(2,1)}$ . The Clebsch-Gordan series for the nontrivial tensor products are:*

	$z_{(2,1),(1,1,1),v}$	$z_{(1,1,1),(2,1),v}$	$z_{(2,1),(2,1),v}$
$v = (3)$	0	0	1
$v = (2, 1)$	1	1	1
$v = (1, 1, 1)$	0	0	1

*The Clebsch-Gordan coefficients for the nontrivial tensor products are given by the following orthogonal matrices:*

$$C_{(2,1) \otimes (1,1,1)} = \begin{bmatrix} 0 & 1 \\ -1 & 0 \end{bmatrix}, \quad C_{(1,1,1) \otimes (2,1)} = \begin{bmatrix} 0 & -1 \\ 1 & 0 \end{bmatrix}, \quad C_{(2,1) \otimes (2,1)} = \frac{\sqrt{2}}{2} \begin{bmatrix} 1 & 0 & -1 & 0 \\ 0 & -1 & 0 & 1 \\ 0 & -1 & 0 & -1 \\ 1 & 0 & 1 & 0 \end{bmatrix}.$$

As in Proposition 10, define:

$$A_{(2,1)\otimes(1,1,1)} = C_{(2,1)\otimes(1,1,1)}^T (\hat{f}_{(2,1)} \otimes \hat{g}_{(1,1,1)}) C_{(2,1)\otimes(1,1,1)}, \quad (13)$$

$$A_{(1,1,1)\otimes(2,1)} = C_{(1,1,1)\otimes(2,1)}^T (\hat{f}_{(1,1,1)} \otimes \hat{g}_{(2,1)}) C_{(1,1,1)\otimes(2,1)}, \quad (14)$$

$$A_{(2,1)\otimes(2,1)} = C_{(2,1)\otimes(2,1)}^T (\hat{f}_{(2,1)} \otimes \hat{g}_{(2,1)}) C_{(2,1)\otimes(2,1)}, \quad (15)$$

Then Proposition 10 gives the following formulas:

$$\widehat{f \cdot g_{\rho(3)}} = \frac{1}{3!} \cdot \left[ \hat{f}_{\rho(3)} \cdot \hat{g}_{\rho(3)} + \hat{f}_{\rho(1,1,1)} \cdot \hat{g}_{\rho(1,1,1)} + 4 \cdot [A_{(2,1)\otimes(2,1)}]_{1,1} \right], \quad (16)$$

$$\begin{aligned} \widehat{f \cdot g_{\rho(2,1)}} = \frac{1}{3!} \cdot & \left[ \hat{f}_{\rho(2,1)} \cdot \hat{g}_{\rho(3)} + \hat{f}_{\rho(3)} \cdot \hat{g}_{\rho(2,1)} + A_{(1,1,1)\otimes(2,1)} \right. \\ & \left. + A_{(2,1)\otimes(1,1,1)} + 2 \cdot [A_{(2,1)\otimes(2,1)}]_{2,3,2;3} \right], \end{aligned} \quad (17)$$

$$\widehat{f \cdot g_{\rho(1,1,1)}} = \frac{1}{3!} \cdot \left[ \hat{f}_{\rho(3)} \cdot \hat{g}_{\rho(1,1,1)} + \hat{f}_{\rho(1,1,1)} \cdot \hat{g}_{\rho(3)} + 4 \cdot [A_{(2,1)\otimes(2,1)}]_{4,4} \right], \quad (18)$$

where the notation  $[A]_{a:b,c:d}$  denotes the block of entries in  $A$  between rows  $a$  and  $b$ , and between columns  $c$  and  $d$  (inclusive).

Using the above formulas, we can continue on Example 8 and compute the last update step in our identity management problem (Figure 2). At the final time step, we observe that Bob is at track 1 with 100% certainty. Our likelihood function is therefore nonzero only for the permutations which map Bob (the second identity) to the first track:

$$L(\sigma) \propto \begin{cases} 1 & \text{if } \sigma = (1, 2) \text{ or } (1, 3, 2) \\ 0 & \text{otherwise} \end{cases}.$$

The Fourier transform of the likelihood function is:

$$\widehat{L}_{\rho(3)} = 2, \quad \widehat{L}_{\rho(2,1)} = \begin{bmatrix} -3/2 & \sqrt{3}/2 \\ -\sqrt{3}/2 & 1/2 \end{bmatrix}, \quad \widehat{L}_{\rho(1,1,1)} = 0. \quad (19)$$

Plugging the Fourier transforms of the prior distribution ( $\widehat{P}^{(2)}$  from Table 5) and likelihood (Equation 19) into Equations 13, 14, 15, we have:

$$A_{(2,1)\otimes(1,1,1)} = \begin{bmatrix} 0 & 0 \\ 0 & 0 \end{bmatrix}, \quad A_{(1,1,1)\otimes(2,1)} = \frac{1}{8} \begin{bmatrix} 1 & \sqrt{3} \\ -\sqrt{3} & -3 \end{bmatrix}, \quad A_{(2,1)\otimes(2,1)} = \frac{1}{32} \begin{bmatrix} -7 & -\sqrt{3} & 11 & 5\sqrt{3} \\ -2\sqrt{3} & -10 & -6\sqrt{3} & -14 \\ 20 & 22\sqrt{3} & -4 & 4\sqrt{3} \\ -11\sqrt{3} & -23 & -\sqrt{3} & -13 \end{bmatrix}$$

To invoke Bayes rule in the Fourier domain, we perform a pointwise product using Equations 16, 17, 18, and normalize by dividing by the trivial coefficient, which yields the Fourier transform of the posterior distribution as:

$$\left[ \widehat{P(\sigma|z)} \right]_{\rho(3)} = 1, \quad \left[ \widehat{P(\sigma|z)} \right]_{\rho(2,1)} = \begin{bmatrix} -1 & 0 \\ 0 & 1 \end{bmatrix}, \quad \left[ \widehat{P(\sigma|z)} \right]_{\rho(1,1,1)} = -1. \quad (20)$$

Finally, we can see that the result is correct by recognizing that the Fourier transform of the posterior (Equation 20) corresponds exactly to the distribution which is 1 at  $\sigma = (1, 2)$  and 0 everywhere else. Bob is therefore at Track 1, Alice at Track 2 and Cathy at Track 3.

$\sigma$	$\varepsilon$	(1, 2)	(2, 3)	(1, 3)	(1, 2, 3)	(1, 3, 2)
$P(\sigma)$	0	1	0	0	0	0

---

**Algorithm 1:** Pseudocode for the Fourier Prediction/Rollup Algorithm.
 

---

PREDICTIONROLLUP  
**input** :  $\hat{Q}_{\rho_\lambda}^{(t)}$  and  $\hat{P}_{\rho_\lambda}^{(t)}$ ,  $\rho_\lambda \in \Lambda$   
**output**:  $\hat{P}_{\rho_\lambda}^{(t+1)}$ ,  $\rho_\lambda \in \Lambda$   
**1** **foreach**  $\rho_\lambda \in \Lambda$  **do**  $\hat{P}_{\rho_\lambda}^{(t+1)} \leftarrow \hat{Q}_{\rho_\lambda}^{(t)} \cdot \hat{P}_{\rho_\lambda}^{(t)}$ ;

---



---

**Algorithm 2:** Pseudocode for the Kronecker Conditioning Algorithm.
 

---

KRONECKERCONDITIONING  
**input** : Fourier coefficients of the likelihood function,  $\hat{L}_{\rho_\lambda}$ ,  $\rho_\lambda \in \Lambda_L$ , and Fourier coefficients of the prior distribution,  $\hat{P}_{\rho_\mu}$ ,  $\rho_\mu \in \Lambda_P$   
**output**: Fourier coefficients of the posterior distribution,  $\widehat{LP}_{\rho_\nu}$ ,  $\rho_\nu \in \Lambda_P$   
**1** **foreach**  $\rho_\nu \in \Lambda_P$  **do**  $\widehat{LP}_{\rho_\nu} \leftarrow \mathbf{0}$  //Initialize Posterior  
 //Pointwise Product  
**2** **foreach**  $\rho_\lambda \in \Lambda_L$  **do**  
**3**     **foreach**  $\rho_\mu \in \Lambda_P$  **do**  
**4**          $z \leftarrow CGseries(\rho_\lambda, \rho_\mu)$ ;  
**5**          $C_{\lambda\mu} \leftarrow CGcoefficients(\rho_\lambda, \rho_\mu)$ ;  $A_{\lambda\mu} \leftarrow C_{\lambda\mu}^T \cdot (\hat{L}_{\rho_\lambda} \otimes \hat{P}_{\rho_\mu}) \cdot C_{\lambda\mu}$ ;  
**6**         **for**  $\rho_\nu \in \Lambda_P$  *such that*  $z_{\lambda\mu\nu} \neq 0$  **do**  
**7**             **for**  $\ell = 1$  **to**  $z_{\lambda\mu\nu}$  **do**  
**8**                  $\left[ \widehat{L^{(t)}P^{(t)}} \right]_{\rho_\nu} \leftarrow \left[ \widehat{L^{(t)}P^{(t)}} \right]_{\rho_\nu} + \frac{d_{\rho_\lambda} d_{\rho_\mu}}{d_{\rho_\nu} n!} A_{\lambda\mu}^{(\nu, \ell)}$ ; //  $A_{\lambda\mu}^{(\nu, \ell)}$  is the  $(\nu, \ell)$  block of  $A_{\lambda\mu}$   
**9**      $\eta \leftarrow \left[ \widehat{L^{(t)}P^{(t)}} \right]_{\rho_{(n)}}^{-1}$ ;  
**10**     **foreach**  $\rho_\nu \in \Lambda$  **do**  $\left[ \widehat{L^{(t)}P^{(t)}} \right]_{\rho_\nu} \leftarrow \eta \left[ \widehat{L^{(t)}P^{(t)}} \right]_{\rho_\nu}$  //Normalization

---

## 7. Approximate Inference by Bandlimiting

We now consider the consequences of performing inference using the Fourier transform at a reduced set of coefficients. Important issues include understanding how error can be introduced into the system, and when our algorithms are expected to perform well as an approximation. Specifically, we fix a bandlimit  $\lambda^{MIN}$  and maintain the Fourier transform of  $P$  only at irreducibles which are at  $\lambda^{MIN}$  or above in the dominance ordering:

$$\Lambda = \{\rho_\lambda : \lambda \succeq \lambda^{MIN}\}.$$

For example, when  $\lambda^{MIN} = (n-2, 1, 1)$ ,  $\Lambda$  is the set  $\{\rho_{(n)}, \rho_{(n-1,1)}, \rho_{(n-2,2)}$ , and  $\rho_{(n-2,1,1)}\}$ , which corresponds to maintaining second-order (ordered) marginal probabilities of the form  $P(\sigma((i, j)) = (k, \ell))$ . During inference, we follow the procedure outlined in the previous section but discard the higher order terms which can be introduced during the conditioning step. Pseudocode for bandlimited prediction/rollup and Kronecker conditioning is given in Algorithms 1 and 2. We note that it is not necessary to maintain the same number of irreducibles for both prior and likelihood during the conditioning step. The first question to ask is: when should one expect a bandlimited approximation to be close to  $P(\sigma)$  as a function? Qualitatively, if a distribution is relatively smooth, then most of its



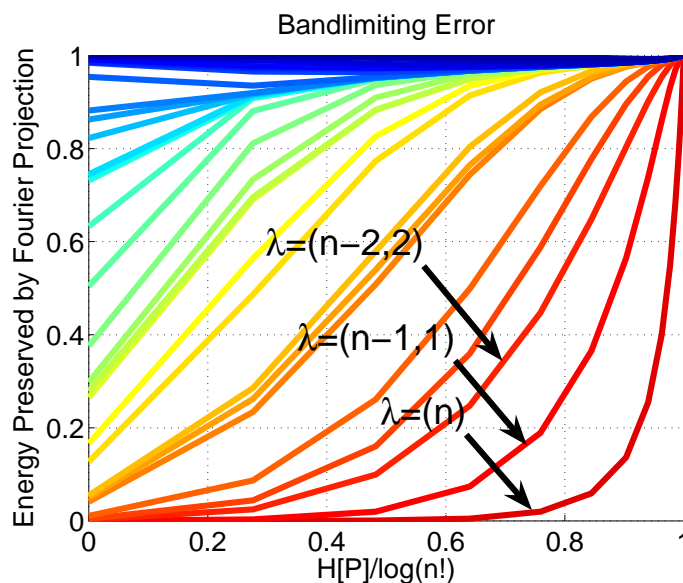


Figure 5: In general, smoother distributions are well approximated by low-order Fourier projections. In this graph, we show the approximation quality of the Fourier projections on distributions with different entropies, starting from sharply peaked delta distributions on the left side of the graph, which get iteratively smoothed until they become the maximum entropy uniform distribution on the right side. On the y-axis, we measure how much energy is preserved in the bandlimited approximation, which we define to be  $\frac{|P'|^2}{|P|^2}$ , where  $P'$  is the bandlimited approximation to  $P$ . Each line represents the approximation quality using a fixed number of Fourier coefficients. At one extreme, we achieve perfect signal reconstruction by using all Fourier coefficients, and at the other, we perform poorly on “spiky” distributions, but well on high-entropy distributions, by storing a single Fourier coefficient.

energy is stored in the low-order Fourier coefficients. However, in a phenomenon quite reminiscent of the Heisenberg uncertainty principle from quantum mechanics, it is exactly when the distribution is sharply concentrated at a small subset of permutations, that the Fourier projection is unable to faithfully approximate the distribution. We illustrate this uncertainty effect in Figure 5 by plotting the accuracy of a bandlimited distribution against the entropy of a distribution.

Even though the bandlimited distribution is sometimes a poor approximation to the true distribution, the marginals maintained by our algorithm are often sufficiently accurate. And so instead of considering the approximation accuracy of the bandlimited Fourier transform to the true joint distribution, we consider the accuracy only at the marginals which are maintained by our method.

### 7.1 Sources of Error During Inference

We now analyze the errors incurred during our inference procedures with respect to the accuracy at maintained marginals. It is immediate that the Fourier domain prediction/rollup operation is *exact* due to its pointwise nature in the Fourier domain. For example, if we have the second order

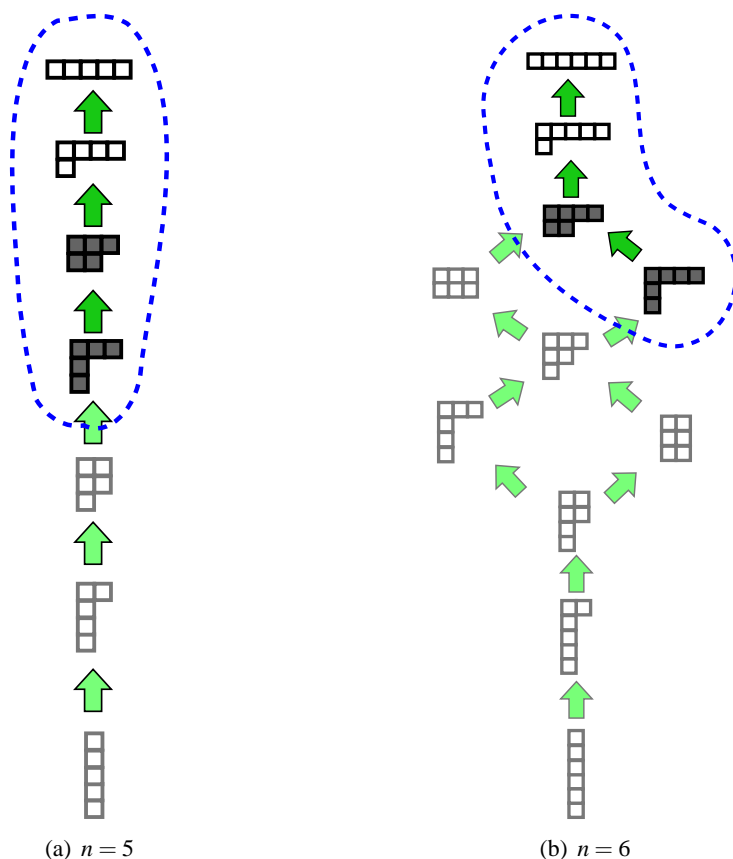


Figure 6: We show the dominance ordering for partitions of  $n = 5$  and  $n = 6$  again. By setting  $\lambda^{MIN} = (3, 1, 1)$  and  $(4, 1, 1)$  respectively, we keep the irreducibles corresponding to the partitions in the dotted regions. If we call Kronecker Conditioning with a first-order observation model, then according to Theorem 12, we can expect to incur some error at the Fourier coefficients corresponding to  $(3, 1, 1)$  and  $(3, 2)$  for  $n = 5$ , and  $(4, 1, 1)$  and  $(4, 2)$  for  $n = 6$  (shown as shaded tableaux), but to be exact at first-order coefficients.

marginals at time  $t = 0$ , then we can find the exact second order marginals at all  $t > 0$  if we only perform prediction/rollup operations. Instead, the errors in inference are only committed by Kronecker conditioning, where they are implicitly introduced at coefficients outside of  $\Lambda$  (by effectively setting the coefficients of the prior and likelihood at irreducibles outside of  $\Lambda$  to be zero), then propagated inside to the irreducibles of  $\Lambda$ .

In practice, we observe that the errors introduced at the low-order irreducibles during inference are small if the prior and likelihood are sufficiently diffuse, which makes sense since the high-frequency Fourier coefficients are small in such cases. We can sometimes show that the update is *exact* at low order irreducibles if we maintain *enough* coefficients.

**Theorem 12** *If  $\lambda^{MIN} = (n - p, \lambda_2, \dots)$ , and the Kronecker conditioning algorithm is called with a likelihood function whose Fourier coefficients are nonzero only at  $\rho_\mu$  when  $\mu \geq (n - q, \mu_2, \dots)$ , then*

the approximate Fourier coefficients of the posterior distribution are exact at the set of irreducibles:

$$\Lambda_{EXACT} = \{\rho_\lambda : \lambda \succeq (n - |p - q|, \dots)\}.$$

**Proof** See Appendix D. ■

For example, if we call Kronecker conditioning by passing in third-order terms of the prior and first-order terms of the likelihood, then all first and second-order (unordered and ordered) marginal probabilities of the posterior distribution can be reconstructed without error.

## 7.2 Projecting to the Marginal Polytope

Despite the encouraging result of Theorem 12, the fact remains that consecutive conditioning steps can propagate errors to all levels of the bandlimited Fourier transform, and in many circumstances, result in a Fourier transform whose “marginal probabilities” correspond to no consistent joint distribution over permutations, and are sometimes negative. To combat this problem, we present a method for projecting to the space of coefficients corresponding to consistent joint distributions (which we will refer to as the *marginal polytope*) during inference.

We begin by discussing the first-order version of the marginal polytope projection problem. Given an  $n \times n$  matrix,  $M$ , of real numbers, how can we decide whether there exists some probability distribution which has  $M$  as its matrix of first-order marginal probabilities? A necessary and sufficient condition, as it turns out, is for  $M$  to be *doubly stochastic*. That is, all entries of  $M$  must be nonnegative and all rows and columns of  $M$  must sum to one (the probability that Alice is at *some track* is 1, and the probability that *some identity* is at Track 3 is 1). The double stochasticity condition comes from the *Birkhoff-von Neumann* theorem (van Lint and Wilson, 2001) which states that a matrix is doubly stochastic *if and only if* it can be written as a convex combination of permutation matrices.

To “renormalize” first-order marginals to be doubly stochastic, some authors (Shin et al., 2003, 2005; Balakrishnan et al., 2004; Helmbold and Warmuth, 2007) have used the *Sinkhorn iteration*, which alternates between normalizing rows and columns independently until convergence is obtained. Convergence is guaranteed under mild conditions and it can be shown that the limit is a nonnegative doubly stochastic matrix which is closest to the original matrix in the sense that the Kullback-Leibler divergence is minimized (Balakrishnan et al., 2004).

There are several problems which cause the Sinkhorn iteration to be an unnatural solution in our setting. First, since the Sinkhorn iteration only works for nonnegative matrices, we would have to first cap entries to lie in the appropriate range,  $[0, 1]$ . More seriously, even though the Sinkhorn iteration would guarantee a doubly stochastic higher order matrix of marginals, there are several natural constraints which are violated when running the Sinkhorn iteration on higher-order marginals. For example, with second-order (ordered) marginals, it seems that we should at least enforce the following symmetry constraint:

$$P(\sigma : \sigma(k, \ell) = (i, j)) = P(\sigma : \sigma(\ell, k) = (j, i)),$$

which says, for example, that the marginal probability that Alice is in Track 1 and Bob is in Track 2 is the same as the marginal probability that Bob is in Track 2 and Alice is in Track 1. Another

natural constraint that can be broken is what we refer to as *low-order marginal consistency*. For example, it should always be the case that:

$$P(j) = \sum_i P(i, j) = \sum_k P(j, k).$$

It should be noted that the doubly stochastic requirement is a special case of lower-order marginal consistency—we require that higher-order marginals be consistent on the  $0^{th}$  order marginal.

While compactly describing the constraints of the marginal polytope exactly remains an open problem, we propose a method for projecting onto a *relaxed* form of the marginal polytope which addresses both symmetry and low-order consistency problems by operating directly on irreducible Fourier coefficients instead of on the matrix of marginal probabilities. After each conditioning step, we apply a ‘correction’ to the approximate posterior  $P^{(t)}$  by finding the bandlimited function in the relaxed marginal polytope which is closest to  $P^{(t)}$  in an  $L_2$  sense. To perform the projection, we employ the Plancherel Theorem (Diaconis, 1988) which relates the  $L_2$  distance between functions on  $S_n$  to a distance metric in the Fourier domain.

**Proposition 13 (Plancherel Theorem)**

$$\sum_{\sigma} (f(\sigma) - g(\sigma))^2 = \frac{1}{|G|} \sum_{\nu} d_{\rho_{\nu}} \text{Tr} \left( (\hat{f}_{\rho_{\nu}} - \hat{g}_{\rho_{\nu}})^T \cdot (\hat{f}_{\rho_{\nu}} - \hat{g}_{\rho_{\nu}}) \right). \tag{21}$$

To find the closest bandlimited function in the relaxed marginal polytope, we formulate a quadratic program whose objective is to minimize the right side of Equation 21, and whose sum is taken only over the set of maintained irreducibles,  $\Lambda$ , subject to the set of constraints which require all marginal probabilities to be nonnegative. We thus refer to our correction step as *Plancherel Projection*. Our quadratic program can be written as:

$$\begin{aligned} & \text{minimize}_{\hat{f}^{proj}} \sum_{\lambda \in \Lambda} d_{\lambda} \text{Tr} \left[ (\hat{f} - \hat{f}^{proj})_{\rho_{\lambda}}^T (\hat{f} - \hat{f}^{proj})_{\rho_{\lambda}} \right] \\ & \text{subject to: } \quad [\hat{f}^{proj}]_{(n)} = 1, \\ & \quad \left[ C_{\lambda_{MIN}} \cdot \left( \bigoplus_{\mu \geq \lambda_{MIN}} \bigoplus_{\ell=1}^{K_{\lambda_{MIN}, \mu}} \hat{f}_{\rho_{\mu}}^{proj} \right) \cdot C_{\lambda_{MIN}}^T \right]_{ij} \geq 0, \quad \text{for all } (i, j), \end{aligned}$$

where  $K_{\lambda_{MIN}}$  and  $C_{\lambda_{MIN}}$  are the precomputed constants from Equation 6. We remark that even though the projection will produce a Fourier transform corresponding to nonnegative marginals which are consistent with each other, there might not necessarily exist a joint probability distribution on  $S_n$  consistent with those marginals except in the special case of first-order marginals.

**Example 11** *In Example 10, we ran the Kronecker conditioning algorithm using all of the Fourier coefficients. If only the first-order coefficients are available, however, then the expressions for zeroth and first order terms of the posterior (Equations 16,17) become:*

$$\begin{aligned} \widehat{f \cdot g}_{\rho_{(3)}} &= \frac{1}{3!} \cdot \left[ \hat{f}_{\rho_{(3)}} \cdot \hat{g}_{\rho_{(3)}} + 4 \cdot [A_{(2,1) \otimes (2,1)}]_{1,1} \right], \\ \widehat{f \cdot g}_{\rho_{(2,1)}} &= \frac{1}{3!} \cdot \left[ \hat{f}_{\rho_{(2,1)}} \cdot \hat{g}_{\rho_{(3)}} + \hat{f}_{\rho_{(3)}} \cdot \hat{g}_{\rho_{(2,1)}} + 2 \cdot [A_{(2,1) \otimes (2,1)}]_{2,3,2;3} \right], \end{aligned}$$

Plugging in the same numerical values from Example 10 and normalizing appropriately yields the approximate Fourier coefficients of the posterior:

$$\left[ \widehat{P(\sigma|z)} \right]_{\rho(3)} = 1 \quad \left[ \widehat{P(\sigma|z)} \right]_{\rho(2,1)} = \begin{bmatrix} -10/9 & -77/400 \\ 77/400 & 4/3 \end{bmatrix},$$

which correspond to the following first-order marginal probabilities:

$$\hat{P}_{\tau(2,1)} \left[ \begin{array}{c|ccc} & A & B & C \\ \hline \text{Track 1} & 0 & 11/9 & -2/9 \\ \text{Track 2} & 1 & 0 & 0 \\ \text{Track 3} & 0 & -2/9 & 11/9 \end{array} \right].$$

In particular, we see that the approximate matrix of ‘marginals’ contains negative numbers. Applying the Plancherel projection step, we obtain the following marginals:

$$\hat{P}_{\tau(2,1)} \left[ \begin{array}{c|ccc} & A & B & C \\ \hline \text{Track 1} & 0 & 1 & 0 \\ \text{Track 2} & 1 & 0 & 0 \\ \text{Track 3} & 0 & 0 & 1 \end{array} \right],$$

which happen to be exactly the true posterior marginals. It should be noted however, that rounding the ‘marginals’ to be in the appropriate range would have worked in this particular example as well.

## 8. Probabilistic Models of Mixing and Observations

While the algorithms presented in the previous sections are general in the sense that they work on all mixing and observation models, it is not always obvious how to compute the Fourier transform of a given model. In this section, we discuss a collection of useful models for which we *can* efficiently compute low-order Fourier coefficients or even provide a closed-form expression. See Table 7 for a summary of the various models covered in this section.

We consider both *mixing* and *observation* models. In multiobject tracking, a mixing model might account for the fact that two tracks may have swapped identities with some probability. Or in card shuffling, a mixing model might reflect that a card has been inserted somewhere into the deck. In multiobject tracking, an observation model might tell us that Alice is at some track with probability one. Or it might reflect the fact that some subset of identities occupies some subset of tracks with no order information, as in the case of the *bluetooth model*. In ranking applications, an observation model might, for example, reflect that some object is ranked higher than, or preferred over some other object.

This section is divided into three parts, each describing a different approach to computing the Fourier coefficients of a model, with some being simpler or more efficient to implement in certain situations than others. In *direct constructions*, we naively apply the definition of the Fourier transform to obtain the Fourier coefficients of some model. In *marginal based constructions*, we first compute the low-order ‘marginals’ of some probabilistic model, then project the result onto the irreducible Fourier basis. Finally, in *coset-based constructions*, we introduce a family of ‘atomic’ indicator functions of subgroups of the form  $S_k \subset S_n$  which are then combined using

Mixing Model	Example Semantics	Relevant Subgroup
Pairwise mixing	Identity confusion at tracks 1 and 2	$S_2$
$k$ -subset mixing	Identity confusion at tracks in $\{1, 2, 4, 6\}$	$S_k$
Insertion mixing	Insert top card somewhere in the deck	n/a
Observation Model	Example Semantics	Relevant Subgroup
Single track observation	Alice is at Track 1	$S_{n-1}$
Multitrack observation	Alice is at Track 1, Bob is at Track 2, etc.	$S_{n-k}$
Bluetooth observation	The girls occupy tracks $\{1, 2, 6, 8\}$	$S_k \times S_{n-k}$
Pairwise ranking observation	Apples are better than oranges	$S_{n-2}$

Table 7: Several useful types of mixing and observation models are summarized in the above table. In many of these cases, computing the appropriate Fourier transform reduces to computing the Fourier transform of the indicator function of some related subgroup of  $S_n$ , and so we also mention the relevant subgroup in the second column. In the third column we provide an example illustrating the semantics of each model.

scale/shift/convolution operations to form more complex models. As we discuss in Section 11, there also remains the open possibility of learning models *directly* in the Fourier domain. For the sake of succinctness, many of the results in this section will be stated without proof.

### 8.1 Direct Construction

In some applications we are fortunate enough to have a model that can be “directly” transformed efficiently using the definition of the Fourier transform (Definition 3). We provide two examples.

#### 8.1.1 PAIRWISE MIXING

The simplest mixing model for identity management assumes that with probability  $p$ , nothing happens, and that with probability  $(1 - p)$ , the identities for tracks  $i$  and  $j$  are swapped. The probability distribution for the *pairwise mixing model* is therefore:

$$Q_{ij}(\pi) = \begin{cases} p & \text{if } \pi = \varepsilon \\ 1 - p & \text{if } \pi = (i, j) \\ 0 & \text{otherwise} \end{cases} . \tag{22}$$

Since  $Q_{ij}$  is such a sparse distribution (in the sense that  $Q_{ij}(\pi) = 0$  for most  $\pi$ ), it is possible to directly compute  $\widehat{Q}_{ij}$  using Definition 3:

$$\left[ \widehat{Q}_{ij} \right]_{\rho_\lambda} = pI + (1 - p)\rho_\lambda((i, j)),$$

where  $I$  refers to the  $d_\lambda \times d_\lambda$  identity matrix (since any representation must map the identity element  $\varepsilon$  to an identity matrix), and  $\rho_\lambda((i, j))$  is the irreducible representation matrix  $\rho_\lambda$  evaluated at the transposition  $(i, j)$  (which can be computed using the algorithms from Appendix C).

8.1.2 INSERTION MIXING

As another example, we can consider the *insertion mixing model* (also called the *top-in shuffle* Diaconis 1988) in which we take the top card in some deck of  $n$  cards, and with uniform probability, insert it *somewhere* in the deck, preserving all other original relative orderings. Insertions can be useful in ranking applications where we might wish to add a new item into consideration without disturbing the marginal probabilities over relative rankings of existing items. The distribution for the insertion mixing model is given by:

$$Q^{insertion}(\pi) = \begin{cases} \frac{1}{n} & \text{if } \pi \text{ is a cycle of the form } (j, j-1, \dots, 1) \text{ for some } j \in \{1, \dots, n\} \\ 0 & \text{otherwise} \end{cases} .$$

Since the insertion mixing model is supported on  $n$  permutations, it is again simple to directly construct the Fourier transform from the definition. We have:

$$\widehat{Q}_{\rho_\lambda}^{insertion} = \frac{1}{n} \sum_{j=1}^n \rho_\lambda(j, j-1, \dots, 1).$$

8.2 Marginal Based Construction

In *marginal based constructions*, we first compute the low-order ‘marginals’<sup>11</sup> of some probabilistic model, then project the result onto the irreducible Fourier basis. Thus given a function  $f : S_n \rightarrow \mathbb{R}$ , we compute, for example, the first-order marginals  $\hat{f}_{\tau_{(n-1,1)}}$ , and conjugate by an intertwining operator (Equation 6) to obtain the Fourier coefficients at  $(n)$  and  $(n-1, 1)$ . Sometimes when the Fourier transform of  $f$  is provably non-zero *only* at low-order terms, a marginal based construction might be the easiest method to obtain Fourier coefficients.

8.2.1 COLOR HISTOGRAM OBSERVATION

The simplest model assumes that we can get observations of the form: ‘track  $\ell$  is color  $k$ ’ (which is essentially the model considered by Kondor et al. 2007). The probability of seeing color  $k$  at track  $\ell$  given data association  $\sigma$  is

$$L(\sigma) = P(z_\ell = k | \sigma) = \alpha_{\sigma^{-1}(\ell), k},$$

where  $\sum_k \alpha_{\sigma^{-1}(\ell), k} = 1$ . For each identity, the likelihood  $L(\sigma) = P(z_\ell = k | \sigma)$  depends, for example, on a histogram over all possible colors. If the number of possible colors is  $K$ , then the likelihood model can be specified by an  $n \times K$  matrix of probabilities. For example,

$$\alpha_{\sigma(\ell), k} = \left[ \begin{array}{c|cccc} & k = \text{Red} & k = \text{Orange} & k = \text{Yellow} & k = \text{Green} \\ \hline \sigma(\text{Alice}) = \ell & 1/2 & 1/4 & 1/4 & 0 \\ \sigma(\text{Bob}) = \ell & 1/4 & 0 & 0 & 3/4 \\ \sigma(\text{Cathy}) = \ell & 0 & 1/2 & 1/2 & 0 \end{array} \right] . \quad (23)$$

Since the observation model only depends on a single identity, the first-order terms of the Fourier transform suffice to describe the likelihood exactly. To compute the first-order Fourier coefficients

11. The word ‘marginals’ is technically appropriate only when the function in question is a legal probability distribution (as opposed to likelihood functions, for example), however we use it to refer to similar summary statistics for general functions.

at irreducibles, we proceed by computing the first-order Fourier coefficients at the first-order permutation representation (the first-order “marginals”), then transforming to irreducible coefficients. The Fourier transform of the likelihood at the first-order permutation representation is given by:

$$\left[\widehat{L}_{\tau_{(n-1,1)}}\right]_{ij} = \sum_{\{\sigma:\sigma(j)=i\}} P(z_\ell = k|\sigma) = \sum_{\{\sigma:\sigma(j)=i\}} \alpha_{\sigma^{-1}(\ell),k}.$$

To compute the  $ij$ -term, there are two cases to consider.

1. If  $i = \ell$  (that is, if Track  $i$  is the same as the track that was observed), then the coefficient  $\widehat{L}_{ij}$  is proportional to the probability that Identity  $j$  is color  $k$ .

$$\widehat{L}_{ij} = \sum_{\{\sigma:\sigma(j)=i\}} \alpha_{j,k} = (n-1)! \cdot \alpha_{j,k}. \tag{24}$$

2. If, on the other hand,  $i \neq \ell$  (Track  $i$  is not the observed track), then the coefficient  $\widehat{L}_{ij}$  is proportional to the sum over

$$\widehat{L}_{ij} = \sum_{\{\sigma:\sigma(j)=i\}} \alpha_{\sigma^{-1}(\ell),k} = \sum_{m \neq j} \sum_{\{\sigma:\sigma(j)=i \text{ and } \sigma(m)=\ell\}} \alpha_{\sigma^{-1}(\ell),k} = \sum_{m \neq j} (n-2)! \cdot \alpha_{m,k}. \tag{25}$$

**Example 12** We will compute the first-order marginals of the likelihood function on  $S_3$  which arises from observing a "Red blob at Track 1". Plugging the values from the “Red” column of the  $\alpha$  matrix (Equation 23) into Equation 24 and 25 yields the following matrix of first-order coefficients (at the  $\tau_{(n-1,1)}$  permutation representation):

$$\left[\widehat{L}_{(n-1,1)}\right]_{ij} = \left[ \begin{array}{c|ccc} & \text{Track 1} & \text{Track 2} & \text{Track 3} \\ \hline \text{Alice} & 1/4 & 1/2 & 3/4 \\ \text{Bob} & 1/4 & 1/2 & 3/4 \\ \text{Cathy} & 1 & 1/2 & 0 \end{array} \right].$$

The corresponding coefficients at the irreducible representations are:

$$\widehat{L}_{(3)} = 1.5, \quad \widehat{L}_{(2,1)} = \begin{bmatrix} 0 & 0 \\ -\sqrt{3}/4 & -3/4 \end{bmatrix}, \quad \widehat{L}_{(1,1,1)} = 0.$$

### 8.2.2 UNORDERED SUBSET (BLUETOOTH) OBSERVATION

We sometimes receive measurements in the form of unordered lists. For example, the *bluetooth model* is the likelihood function that arises if tracks  $\{1, \dots, k\}$  are within range of a bluetooth detector and we receive a measurement that identities  $\{1, \dots, k\}$  are in range. In sports, we might observe that the first  $k$  tracks belong to the red team and that the last  $n - k$  tracks belong to the blue team. And finally, in *approval voting*, one specifies a subset of approved candidates rather than, for example, picking a single favorite.

We consider two options for bluetooth-type situations. In the first option, we allow for some error-tolerance by setting the likelihood to be proportional to the number of tracks that are correctly returned in the measurement:

$$P^{bluetooth}(z_{\{t_1, \dots, t_k\}} = \{i_1, \dots, i_k\}|\sigma) \propto |\{t_1, \dots, t_k\} \cap \sigma(\{i_1, \dots, i_k\})| + U(\sigma), \tag{26}$$



where  $U(\sigma)$  is a constant function on  $S_n$  allowing for noisy observations. Our first bluetooth model can be expressed using only first order terms (intuitively because each track makes a linear contribution) and thus  $\hat{P}_\lambda^{\text{bluetooth}}$  is nonzero only at the first two partitions  $\lambda = (n), (n-1, 1)$ . For simplicity, we consider the Fourier transform of the function:  $f(\sigma) = |\sigma(\{1, \dots, k\}) \cap \{1, \dots, k\}|$ . The first-order ‘marginals’ of  $f$  are covered in the following four cases:

- ( $j \leq k$  and  $i \leq k$ ):  $L_{ij} = \sum_{\sigma: \sigma(j)=i} f(\sigma) = (k-1)^2(n-2)! + (n-1)!$
- ( $j \leq k$  and  $i > k$ ):  $L_{ij} = \sum_{\sigma: \sigma(j)=i} f(\sigma) = k(k-1)(n-2)!$
- ( $j > k$  and  $i \leq k$ ):  $L_{ij} = \sum_{\sigma: \sigma(j)=i} f(\sigma) = k(k-1)(n-2)!$
- ( $j > k$  and  $i > k$ ):  $L_{ij} = \sum_{\sigma: \sigma(j)=i} f(\sigma) = k^2(n-2)!$

We discuss the second bluetooth-type model after discussing coset based constructions.

### 8.3 Coset-Based Construction

Most of the time, realistic models are not supported on only a handful of permutations. The approach we take now is to use a collection of ‘primitive’ functions to form more interesting models via scale/shift/convolution operations. In particular, we will make use of indicator functions of subsets of the form  $S_{X,Y} \subset S_n$ , where  $X = (x_1, \dots, x_k)$  and  $Y = (y_1, \dots, y_k)$  are ordered  $k$ -tuples with  $\{x_1, \dots, x_k\} \subset \{1, \dots, n\}$ ,  $\{y_1, \dots, y_k\} \subset \{1, \dots, n\}$  and no repetitions are allowed.  $S_{X,Y}$  denotes the set of elements in  $S_n$  which are constrained to map each  $x_i$  to  $y_i$ :

$$S_{X,Y} \equiv \{\sigma \in S_n : \sigma(x_i) = y_i, \text{ for each } i=1, \dots, k\}. \quad (27)$$

The  $S_{X,Y}$  can also be thought of as two-sided cosets associated with subgroups of the form  $S_{n-k} \subset S_n$ . For example, if  $X = (1, 2)$  and  $Y = (3, 4)$  with  $n = 4$ , then  $S_{X,Y}$  is simply the set of all permutations that map  $1 \mapsto 3$  and  $2 \mapsto 4$ . Thus,  $S_{X,Y} = \{(1, 3)(2, 4), (1, 3, 2, 4)\}$ . Since  $|X| = |Y| = k$ , then  $|S_{X,Y}| = (n-k)!$ , and in the special case that  $X = Y$ , we have that  $S_{X,Y}$  is in fact a subgroup isomorphic to  $S_{n-k}$ .

As we show in Appendix C, the Fourier transform of the indicator  $\delta_{S_{X,Y}}$  takes a particularly simple (and low rank) form and can be efficiently computed. The method described in Appendix C is based on the FFT and exploits the same structure of the symmetric group that is used by Kondor et al. (2007). It is thus possible to understand why some observation models afford faster conditioning updates based on sparsity in Fourier domain.

The functions  $\delta_{S_{X,Y}}$  can be viewed as a set of function primitives for constructing more complicated models via shift/scale/convolution operations in the Fourier domain. We now discuss the remaining models in Table 7 with the assumption that there exists some blackbox function which constructs the Fourier coefficients of the indicator function of (two-sided) cosets of the form  $S_{X,Y} \subset S_n$  (see Algorithm 5 in Appendix D).

#### 8.3.1 $k$ -SUBSET MIXING

It is not always appropriate to mix only two people at once (as in Equation 22) and so we would like to formulate a mixing model which occurs over a subset of tracks,  $X = \{t_1, \dots, t_k\} \subset \{1, \dots, n\}$ . One way to ‘mimic’ the desired effect is to repeatedly draw pairs  $(i, j)$  from  $\{t_1, \dots, t_k\}$  and to convolve against the pairwise mixing models  $Q_{ij}$ . A better alternative is to directly construct the Fourier coefficient matrices for the  $k$ -subset mixing model, in which we allow the tracks in  $X$  to be

randomly permuted with uniform probability. In the following,  $\bar{X}$  denotes some fixed ordering of the complement of  $X$ . For example, if  $n = 5$ , with  $X = \{1, 2, 4\}$ , then  $\bar{X}$  is either  $(3, 5)$  or  $(5, 3)$ . The  $k$ -subset mixing model is defined as:

$$Q_X(\pi) = \begin{cases} \frac{1}{k!} & \text{if } \pi \in S_{\bar{X}, \bar{X}} \subset S_n \\ 0 & \text{otherwise} \end{cases}. \quad (28)$$

Note that  $S_{\bar{X}, \bar{X}}$  is isomorphic to  $S_k$  and that the pairwise mixing model is the special case where  $k = 2$ . Intuitively, Equation 28 fixes all of the tracks outside of  $X$  and says that with uniform probability, the set of tracks in  $X$  experience some permutation of their respective identities. Equation 28 can also be written as  $Q_X(\pi) = \frac{1}{k!} \delta_{S_{\bar{X}, \bar{X}}}(\pi)$ , and thus the mixing model is simply a multiple of the indicator function of  $S_{\bar{X}, \bar{X}}$ .

### 8.3.2 SINGLE/MULTI-TRACK OBSERVATION

In the *single track observation model* (used in Shin et al. 2005, Schumitsch et al. 2005 and Kondor et al. 2007, for example), we acquire an identity measurement  $z_j$  at track  $j$ . In the simplest version of the model, we write the likelihood function as:

$$P(z_i = j | \sigma) = \begin{cases} \pi & \text{if } \sigma(j) = i \\ \frac{1-\pi}{n-1} & \text{otherwise} \end{cases}, \quad (29)$$

where  $j$  ranges over all  $n$  possible identities.  $P(z_i | \sigma)$  can also be written as a weighted sum of a uniform distribution  $U$ , and an indicator function:

$$P(z_i = j | \sigma) = \left( \frac{\pi n - 1}{n - 1} \right) \delta_{S_{j,i}}(\sigma) + \left( \frac{1 - \pi}{n - 1} \right) U(\sigma).$$

Equation 29 is useful when we receive measurements directly as single identities (“Alice is at Track 1 with such and such probability”). It is, however, far more common to receive lower level measurements that *depend* only upon a single identity, which we formalize with the following conditional independence assumption:

$$P(z_i | \sigma) = P(z_i | \sigma(j)).$$

For example, as in Equation 23, we might have a color histogram over each individual (“Alice loves to wear green”) and observe a single color per timestep. Or we might acquire observations in the form of color histograms and choose to model a distribution over all possible color histograms. If for each identity  $j$ ,  $P(z_i | \sigma(j) = i) = \alpha_j$ , then we can write the likelihood function as a weighted linear combination of  $n$  indicators,

$$L(\sigma) = P(z_i | \sigma) = \sum_j \alpha_j \delta_{S_{j,i}}(\sigma),$$

and by the linearity of the Fourier transform, we can obtain the Fourier coefficients of  $L$ :

$$\hat{L}_\lambda = \sum_j \alpha_j \left[ \widehat{\delta_{S_{j,i}}} \right]_\lambda.$$

Finally, the single-track observations can be generalized to handle joint observations of multiple tracks at once with a higher-order model:

$$P(z_{(t_1, \dots, t_k)} = (i_1, \dots, i_k) | \sigma) = \begin{cases} \pi & \text{if } \sigma(i_\ell) = t_\ell \text{ for each } \ell \in \{1, \dots, k\} \\ \frac{1-\pi}{\binom{n}{n-k} - 1} & \text{otherwise} \end{cases}. \quad (30)$$

Unsurprisingly, while the Fourier coefficients of Equation 29 can be expressed exactly using first-order terms, the Fourier coefficients of the multi-track observation model, Equation 30, requires  $k$ th-order terms. It is important to note that joint multi-track observations are distinct from making  $k$  independent identity observations at the same timestep—we can handle the latter case by calling the Kronecker conditioning algorithm with a single-track observation model  $k$  times. Depending upon the specific sensor setup, one model may be more natural than the other.

### 8.3.3 BLUETOOTH OBSERVATION

In contrast with the first bluetooth model (Equation 26), our second bluetooth-type model handles a higher-order form of measurement. Like the single/multi-track observation models, it says that with some probability we receive the correct unordered list, and with some probability, we receive some other list drawn uniformly at random:

$$P^{bluetooth2}(z_{\{t_1, \dots, t_k\}} = \{i_1, \dots, i_k\} | \sigma) = \begin{cases} \pi & \text{if } \sigma(\{i_1, \dots, i_k\}) = \{t_1, \dots, t_k\} \\ \frac{1-\pi}{\binom{n}{k} - 1} & \text{otherwise} \end{cases}.$$

As with the single/multi-track observation models, the bluetooth model can be written as a weighted linear combination of a uniform distribution and the indicator function of an  $S_k \times S_{n-k}$ -coset, where:

$$S_k \times S_{n-k} = \{\sigma \in S_n : \sigma(\{1, \dots, k\}) = \{1, \dots, k\}\}.$$

To compute the Fourier transform of  $P^{bluetooth2}$ , it is enough to note that the indicator function of  $S_k \times S_{n-k}$  can be thought of as a convolution of indicator functions of  $S_k$  and  $S_{n-k}$  in a certain sense. More precisely.

**Proposition 14** *Let  $X = (1, \dots, k)$  and  $Y = (k+1, \dots, n)$ . Then:  $\delta_{S_k \times S_{n-k}} = \delta_{S_{X,X}} * \delta_{S_{Y,Y}}$ .*

Invoking the convolution theorem (Proposition 8) shows that the Fourier coefficient matrices of  $\delta_{S_k \times S_{n-k}}$  can be constructed by first computing the Fourier coefficients of  $S_{X,X}$  and  $S_{Y,Y}$ , and point-wise multiplying corresponding coefficient matrices. We have:

$$\left[ \widehat{\delta}_{S_k \times S_{n-k}} \right]_\lambda = \left[ \widehat{\delta}_{S_{X,X}} \right]_\lambda \cdot \left[ \widehat{\delta}_{S_{Y,Y}} \right]_\lambda, \text{ for all partitions } \lambda.$$

An interesting fact about the bluetooth model is that its Fourier terms are zero at all partitions with more than two rows.

**Proposition 15** *Without loss of generality, assume that  $k \leq \frac{n}{2}$ . The Fourier transform of the bluetooth model,  $\widehat{P}_\lambda^{bluetooth2}$  is nonzero only at partitions of the form  $(n-s, s)$  where  $s \leq k$ .*

### 8.3.4 PAIRWISE RANKING OBSERVATION

Finally in the *pairwise ranking model*, we consider observations of the form “object  $j$  is ranked higher than object  $i$ ” which can appear in various forms of voting and preference elicitation (“I like candidate  $x$  better than candidate  $y$ ”) or webpage/advertisement ranking. Here we think of  $\sigma$  as a mapping from objects to ranks. Our pairwise ranking model simply assigns higher probability to observations which agree with the ordering of  $i$  and  $j$  in  $\sigma$ .

$$P^{rank}(z_{k\ell}|\sigma) = \begin{cases} \pi & \text{if } \sigma(k) < \sigma(\ell) \\ 1 - \pi & \text{otherwise} \end{cases}.$$

When  $k = n - 1$ ,  $\ell = n$  and  $\pi = 1$ , we have:

$$\begin{aligned} P^{rank}(z_{k\ell}|\sigma) &= \begin{cases} 1 & \text{if } \sigma(n - 1) < \sigma(n) \\ 0 & \text{otherwise} \end{cases} \\ &= \sum_{i < j} \delta_{S_{(n-1,n),(i,j)}}(\sigma). \end{aligned}$$

Perhaps unsurprisingly, pairwise ranking models can be sufficiently captured by first-order and second-order (ordered) Fourier coefficients<sup>12</sup>.

**Proposition 16** *The Fourier coefficients of the pairwise ranking model,  $\hat{P}_\lambda^{rank}$ , are nonzero only at three partitions:  $\lambda = (n)$ ,  $(n - 1, 1)$ , and  $(n - 2, 1, 1)$ .*

## 9. Related Work

Rankings and permutations have recently become an active area of research in machine learning due to their importance in information retrieval and preference elicitation. Rather than considering full distributions over permutations, many approaches, like RankSVM (Joachims, 2002) and RankBoost (Freund et al., 2003), have instead focused on learning a single ‘optimal’ ranking with respect to some objective function.

There are also several authors (from both the statistics and machine learning communities) who have studied distributions over permutations/rankings (Mallows, 1957; Critchlow, 1985; Fligner and Verducci, 1986; Meila et al., 2007; Taylor et al., 2008; Lebanon and Mao, 2008). Taylor et al. (2008) consider distributions over  $S_n$  which are induced by the rankings of  $n$  independent draws from  $n$  individually centered Gaussian distributions with equal variance. They compactly summarize their distributions using an  $O(n^2)$  matrix which is conceptually similar to our first-order summaries and apply their techniques to ranking web documents. Most other previous approaches at directly modeling distributions on  $S_n$ , however, have relied on distance based exponential family models. For example, the Mallows model (Mallows, 1957) defines a Gaussian-like distribution over permutations as:

$$P(\sigma; c, \sigma_0) \propto \exp(-cd(\sigma, \sigma_0)),$$

where the function  $d(\sigma, \sigma_0)$  is the *Kendall’s tau distance* which counts the number of adjacent swaps that are required to bring  $\sigma^{-1}$  to  $\sigma_0^{-1}$ .

---

12. Additionally,  $\hat{P}_{(n-1,1)}^{rank}$  and  $\hat{P}_{(n-2,1,1)}^{rank}$  are known to be rank 1 matrices, a fact which can potentially be exploited for faster conditioning updates in practice.

Distance based exponential family models have the advantage that they can compactly represent distributions for very large  $n$ , and admit conjugate prior distributions (Meila et al., 2007). Estimating parameters has been a popular problem for statisticians—recovering the optimal  $\sigma_0$  from data is known as the *consensus ranking* or *rank aggregation* problem and is known to be *NP*-hard (Bartholdi et al., 1989). Many authors have focused on approximation algorithms instead.

Like Gaussian distributions, distance based models also tend to lack flexibility, and so Lebanon and Mao (2008) propose a nonparametric model of ranked (and partially ranked) data based on placing weighted Mallows kernels on top of training examples, which, as they show, can realize a far richer class of distributions, and can be learned efficiently. However, they do not address the inference problem, and it is not clear if one can efficiently perform inference operations like marginalization and conditioning in such models.

As we have shown in this paper, Fourier based methods (Diaconis, 1988; Kondor et al., 2007; Huang et al., 2007) offer a principled alternative method for compactly representing distributions over permutations and performing efficient probabilistic inference operations. Our work draws from two strands of research—one from the data association/identity management literature, and one from a more theoretical area on Fourier analysis in statistics. In the following, we review several of the works which have led up to our current Fourier based approach.

## 9.1 Previous Work in Identity Management

The identity management problem has been addressed in a number of previous works, and is closely related to, but not identical with, the classical data association problem of maintaining correspondences between tracks and observations. Both problems need to address the fundamental combinatorial challenge that there is a factorial or exponential number of associations to maintain between tracks and identities, or between tracks and observations respectively. A vast literature already exists on the the data association problem, beginning with the *multiple hypothesis testing* approach (MHT) of Reid (1979). The MHT is a ‘deferred logic’ method in which past observations are exploited in forming new hypotheses when a new set of observations arises. Since the number of hypotheses can grow exponentially over time, various heuristics have been proposed to help cope with the complexity blowup. For example, one can choose to maintain only the  $k$  *best* hypotheses for some parameter  $k$  (Cox and Hingorani, 1994), using Murty’s algorithm (Murty, 1968). But for such an approximation to be effective,  $k$  may still need to scale exponentially in the number of objects. A slightly more recent filtering approach is the *joint probabilistic data association filter* (JPDA) (Bar-Shalom and Fortmann, 1988), which is a suboptimal single-stage approximation of the optimal Bayesian filter. JPDA makes associations sequentially and is unable to correct erroneous associations made in the past (Poore, 1995). Even though the JPDA is more efficient than the MHT, the calculation of the JPDA association probabilities is still a  $\#P$ -complete problem (Collins and Uhlmann, 1992), since it effectively must compute matrix permanents. Polynomial approximation algorithms to the JPDA association probabilities have recently been studied using Markov chain Monte Carlo (MCMC) methods (Oh et al., 2004; Oh and Sastry, 2005).

The identity management problem was first explicitly introduced in Shin et al. (2003). Identity management differs from the classical data association problem in that its observation model is not concerned with the low-level tracking details but instead with high level information about object identities. Shin et al. (2003) introduced the notion of the *belief matrix* approximation of the association probabilities, which collapses a distribution over all possible associations to just

its first-order marginals. In the case of  $n$  tracks and  $n$  identities, the belief matrix  $B$  is an  $n \times n$  doubly-stochastic matrix of non-negative entries  $b_{ij}$ , where  $b_{ij}$  is the probability that identity  $i$  is associated with track  $j$ . As we already saw in Section 4, the belief matrix approximation is equivalent to maintaining the zeroth- and first-order Fourier coefficients. Thus our current work is a strict generalization and extension of those previous results.

An alternative representation that has also been considered is an information theoretic approach (Shin et al., 2005; Schumitsch et al., 2005, 2006) in which the density is parameterized as:

$$P(\sigma; \Omega) \propto \exp \text{Tr} (\Omega^T \cdot \tau_{(n-1,1)}(\sigma)).$$

In our framework, the information form approach can be viewed as a method for maintaining the Fourier transform of the *log* probability distribution at only the first two irreducibles. The information matrix approach is especially attractive in a distributed sensor network setting, since, if the columns of the information matrix are distributed to leader nodes tracking the respective targets, then the observation events become entirely local operations, avoiding the more expensive Kronecker conditioning algorithm in our setting. On the other hand, the information matrix coefficients do not have the same intuitive marginals interpretation afforded in our setting, and moreover, prediction/rollup steps cannot be performed analytically in the information matrix form. As in many classical data structures problems there are representation trade-off issues: some operations are less expensive in one representation and some operations in the the other. The best choice in any particular scenario will depend on the ratio between observation and mixing events.

## 9.2 Previous Work on Fourier-Based Approximations

The concept of using Fourier transforms to study probability distributions on groups is not new, with the earliest papers in this area having been published in the 1960s (Grenander, 1963). Willsky (1978) was the first to formulate the exact filtering problem in the Fourier domain for finite and locally compact Lie groups and contributed the first noncommutative Fast Fourier Transform algorithm (for Metacyclic groups). However, he does not address approximate inference, suggesting instead to always transform to the appropriate domain for which either the prediction/rollup or conditioning operations can be accomplished using a pointwise product. While providing significant improvements in complexity for smaller groups, his approach is still infeasible for our problem given the factorial order of the Symmetric group.

Diaconis (1988) used the Fourier transform to analyze probability distributions on the Symmetric group in order to study card shuffling and ranking problems. His work laid the ground for much of the progress made over the last two decades on probabilistic group theory and noncommutative FFT algorithms (Clausen and Baum, 1993; Rockmore, 2000).

Kondor et al. (2007) was the first to show that the data association problem could be efficiently approximated using FFT factorizations. In contrast to our framework where every model is assumed to be have been specified in the Fourier domain, they work with an observation model which can be written as the indicator function of cosets of subgroups of the form  $S_k \subset S_n$ .

Conceptually, one might imagine formulating a conditioning algorithm which applies the Inverse Fast Fourier Transform (IFFT) to the prior distribution, conditions in the primal domain using pointwise multiplication, then transforms back up to the Fourier domain using the FFT to obtain posterior Fourier coefficients. While such a procedure would ordinarily be intractable because of the factorial number of permutations, Kondor et al. (2007) elegantly shows that for certain coset-

based observation models, it is not necessary to perform the full FFT recursion to do a pointwise product. They exploit this observation to formulate an efficient conditioning algorithm whose running time depends on the complexity of the observation model (which can roughly be measured by the number of irreducibles required to fully specify it).

Our work generalizes the conditioning formulation from Kondor et al. (2007) in the sense that it can work for *any* observation model and extends easily to similar filtering problems over any finite group. In the case that the observation model is specified at sufficiently many irreducibles, our conditioning algorithm (prior to the projection step) returns the same approximate probabilities as the FFT-based algorithm. For example, we can show that the observation model given in Equation 29 is fully specified by two Fourier components, and that both algorithms have identical output. Additionally, Kondor et al. (2007) do not address the issue of projecting onto legal distributions, which, as we show in our experimental results is fundamental in practice.

## 10. Experimental Results

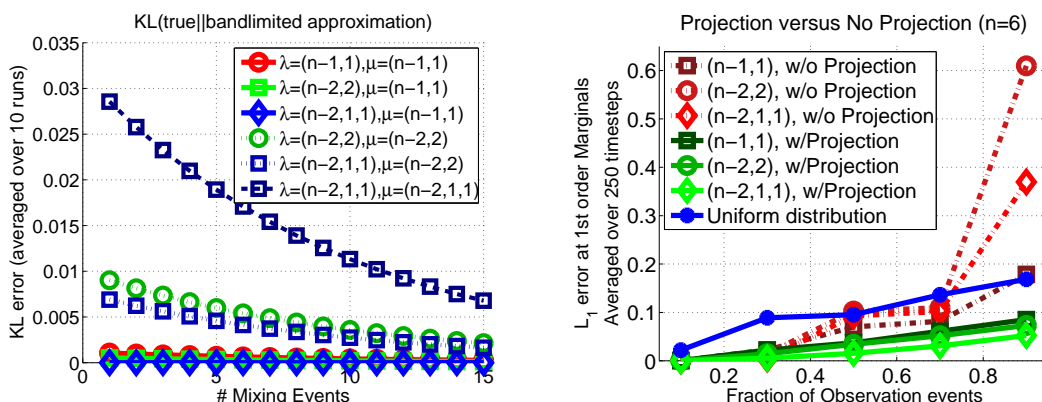
In this section we present the results of several experiments to validate our algorithm. We evaluate performance first by measuring the quality of our approximation for problems where the true distribution is known. Instead of measuring a distance between the true distribution and the inverse Fourier transform of our approximation, it makes more sense in our setting to measure error only at the marginals which are maintained by our approximation. In the results reported below, we measure the  $L_1$  error between the true matrix of marginals and the approximation. If nonnegative marginal probabilities are guaranteed, it also makes sense to measure KL-divergence.

### 10.1 Simulated Data

We first tested the accuracy of a single Kronecker conditioning step by calling some number of pairwise mixing events (which can be thought roughly as a measure of entropy), followed by a single first-order observation. In the y-axis of Figure 7(a), we plot the Kullback-Leibler divergence between the true first-order marginals and approximate first-order marginals returned by Kronecker conditioning. We compared the results of maintaining first-order, and second-order (unordered and ordered) marginals. As shown in Figure 7(a), Kronecker conditioning is more accurate when the prior is smooth and unsurprisingly, when we allow for higher order Fourier terms. As guaranteed by Theorem 12, we also see that the first-order terms of the posterior are exact when we maintain second-order (ordered) marginals.

To understand how our algorithms perform over many timesteps (where errors can propagate to all Fourier terms), we compared to exact inference on synthetic data sets in which tracks are drawn at random to be observed or swapped. As a baseline, we show the accuracy of a uniform distribution. We observe that the Fourier approximation is better when there are either more mixing events (the fraction of conditioning events is smaller), or when more Fourier coefficients are maintained, as shown in Figure 7(b). We also see that the Plancherel Projection step is fundamental, especially when mixing events are rare.

Figures 10(a) and 10(b) show the per-timeslice accuracy of two typical runs of the algorithm. The fraction of conditioning events is 50% in Figure 10(a), and 70% in Figure 10(b). What we typically observe is that while the projected and nonprojected accuracies are often quite similar, the nonprojected marginals can perform significantly worse during certain segments.



(a) Kronecker Conditioning Accuracy—we measure the accuracy of a single Kronecker conditioning operation after some number of mixing events. (b) HMM Accuracy—we measure the average accuracy of posterior marginals over 250 timesteps, varying the proportion of mixing and observation events

Figure 7: Simulation results.

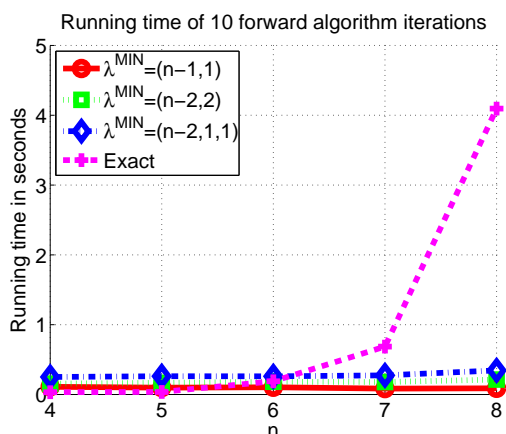


Figure 8: Running times: We compared running times of our polynomial time bandlimited inference algorithms against an exact algorithm with  $O(n^3n!)$  time complexity

Finally, we compared running times against an exact inference algorithm which performs prediction/rollup in the Fourier domain and conditioning in the primal domain. While the prediction/rollup step for pairwise mixing models can be implemented in  $O(n!)$  time (linear in the size of the symmetric group), we show running times for the more general mixing models. Instead of the naive  $O((n!)^2)$  complexity, its running time is a more efficient  $O(n^3n!)$  due to the Fast Fourier Transform (Clausen and Baum, 1993). It is clear that our algorithm scales gracefully compared to the exact solution (Figure 10.1), and in fact, we could not run exact inference for  $n > 8$  due to memory constraints. In Figure 10.1, we show empirically that the Clebsch-Gordan coefficients are indeed sparse, supporting a faster conjectured runtime.



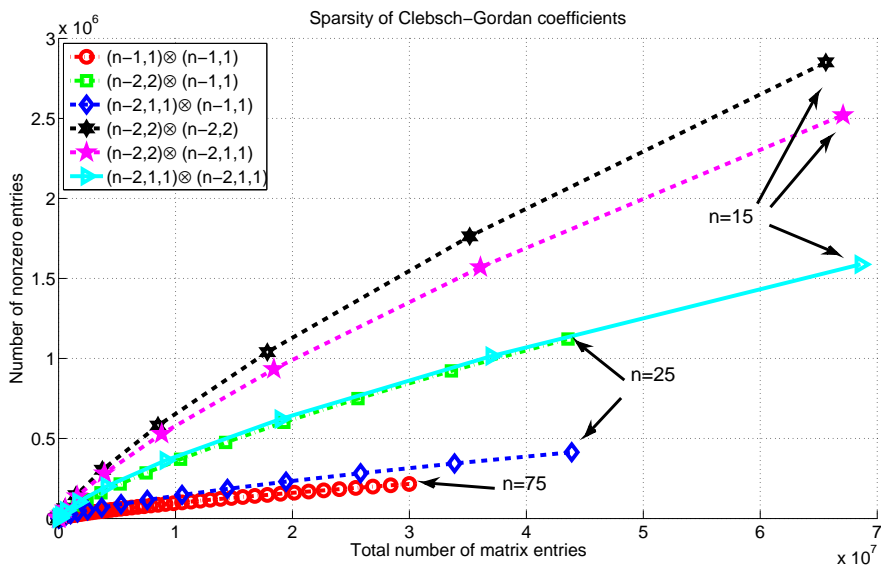


Figure 9: Clebsch-Gordan Sparsity: We measured the sparsity of the Clebsch-Gordan coefficients matrices by plotting the number of nonzero coefficients in a Clebsch-Gordan coefficient matrix against the number of total entries in the matrix for various  $n$  and pairs of irreducibles. For each fixed tensor product pair, we see that the number of nonzero entries scales sublinearly with respect to the total number of matrix elements.

## 10.2 Real Camera Network

We also evaluated our algorithm on data taken from a real network of eight cameras (Fig. 11(a)). In the data, there are  $n = 11$  people walking around a room in fairly close proximity. To handle the fact that people can freely leave and enter the room, we maintain a list of the tracks which are external to the room. Each time a new track leaves the room, it is added to the list and a mixing event is called to allow for  $m^2$  pairwise swaps amongst the  $m$  external tracks.

The number of mixing events is approximately the same as the number of observations. For each observation, the network returns a color histogram of the blob associated with one track. The task after conditioning on each observation is to predict identities for all tracks which are inside the room, and the evaluation metric is the fraction of accurate predictions. We compared against a baseline approach of predicting the identity of a track based on the most recently observed histogram at that track. This approach is expected to be accurate when there are many observations and discriminative appearance models, neither of which our problem afforded. As Figure 11(b) shows, both the baseline and first order model (without projection) fared poorly, while the projection step dramatically boosted the prediction accuracy for this problem. To illustrate the difficulty of predicting based on appearance alone, the rightmost bar reflects the performance of an *omniscient* tracker who knows the result of each mixing event and is therefore left only with the task of distinguishing between appearances. We conjecture that the performance of our algorithm (with projection) is near optimal.

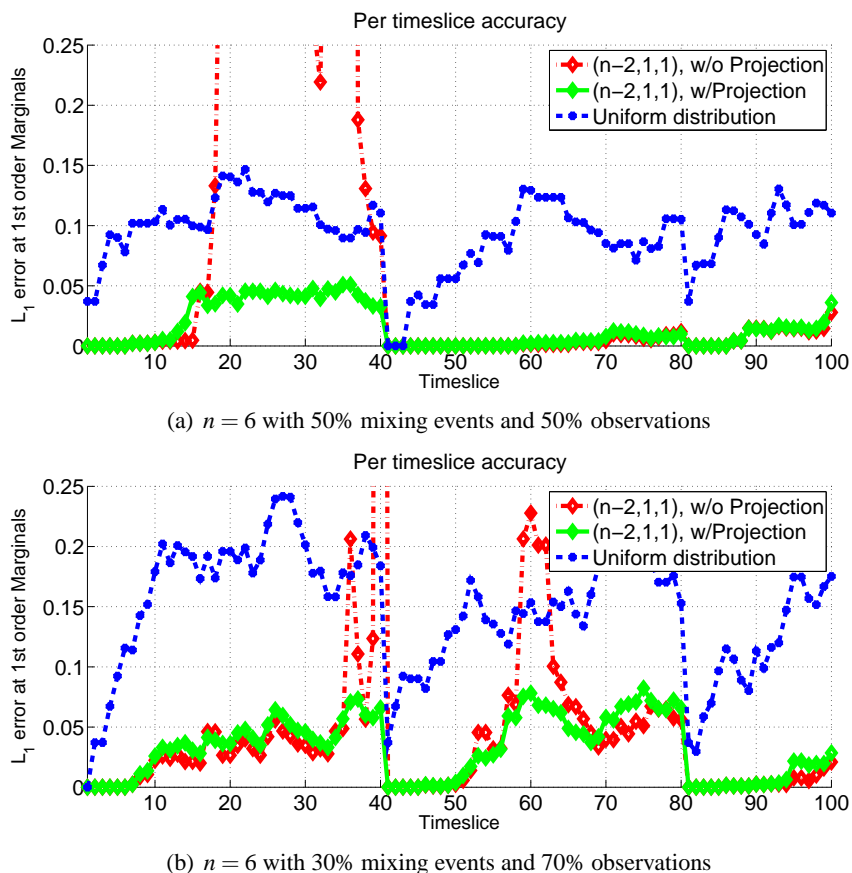


Figure 10: Accuracy as a function of time on two typical runs.

## 11. Future Research

There remain several possible extensions to the current work stemming from both practical and theoretical considerations. We list a few open questions and extensions in the following.

### 11.1 Adaptive Filtering

While our current algorithms easily beat exact inference in terms of running time, they are still limited by a relatively high (though polynomial) time complexity. In practice however, it seems reasonable to believe that the “difficult” identity management problems typically involve only a small subset of people at a time. A useful extension of our work would be to devise an *adaptive* version of the algorithm which allocates more Fourier coefficients towards the identities which require higher order reasoning. We believe that this kind of extension would be the appropriate way to scale our algorithm to handling massive numbers of objects at a time.

### 11.2 Characterizing the Marginal Polytope

In our paper, we presented a projection of the bandlimited distribution to a certain polytope, which is exactly the marginal polytope for first-order bandlimited distributions, but strictly an outer bound

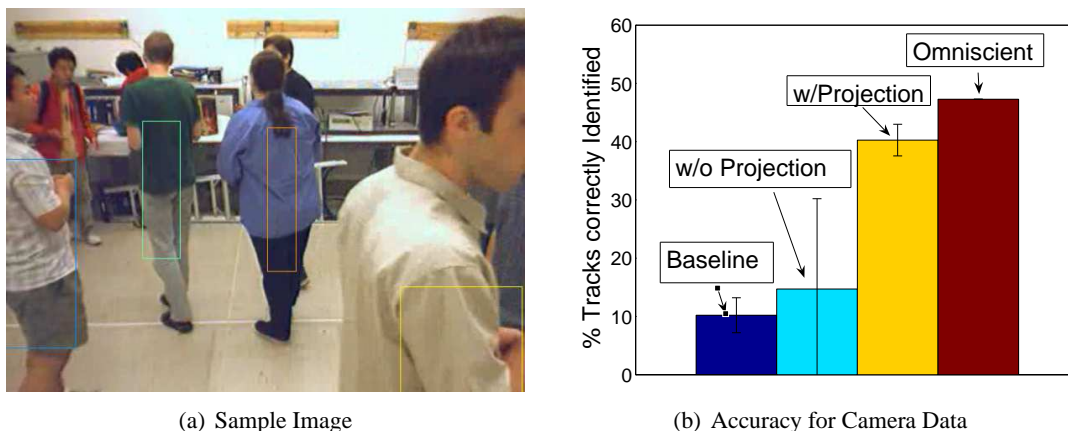


Figure 11: Evaluation on data set from a real camera network. In this experiment, there are  $n = 11$  people walking in a room begin tracked by 8 cameras.

for higher orders. An interesting project would be to generalize the Birkhoff-von Neumann theorem by exactly characterizing the marginal polytope at higher order marginals. We conjecture that the marginal polytope for low order marginals can be described with polynomially many constraints.

### 11.3 Learning in the Fourier Domain

Another interesting problem is whether we can learn bandlimited mixing and observation models *directly in the Fourier domain*. Given fully observed permutations  $\sigma_1, \dots, \sigma_m$ , drawn from a distribution  $P(\sigma)$ , a naive method for estimating  $\hat{P}_\rho$  at low-order  $\rho$  is to simply observe that:

$$\hat{P}_\rho = \mathbb{E}_{\sigma \sim P}[\rho(\sigma)],$$

and so one can estimate the Fourier transform by simply averaging  $\rho(\sigma_i)$  over all  $\sigma_i$ . However, since we typically do not observe full permutations in real applications like ranking or identity management, it would be interesting to estimate Fourier transforms using partially observed data. In the case of Bayesian learning, it may be possible to apply some of the techniques discussed in this paper.

### 11.4 Probabilistic Inference on Other Groups

The Fourier theoretic framework presented in this paper is not specific to the symmetric group - in fact, the prediction/rollup and conditioning formulations, as well as most of the results from Appendix D hold over any finite or compact Lie group. As an example, the noncommutative group of rotation operators in three dimensions,  $SO(3)$ , appears in settings which model the pose of a three dimensional object. Elements in  $SO(3)$  might be used to represent the pose of a robot arm in robotics, or the orientation of a mesh in computer graphics; In many settings, it would be useful to have a compact representation of uncertainty over poses. We believe that there are many other application domains with algebraic structure where similar probabilistic inference algorithms might apply, and in particular, that noncommutative settings offer a particularly challenging but exciting opportunity for machine learning research.

## 12. Conclusions

In this paper, we have presented a Fourier theoretic framework for compactly summarizing distributions over permutations. We showed that common probabilistic inference operations can be performed completely in the Fourier domain and that, using the low-order terms of the Fourier expansion of a distribution, one can obtain polynomial time inference algorithms. Fourier theoretic summaries are attractive because they have tuneable approximation quality, have intuitive interpretations in terms of low-order marginals, and have allowed us to leverage results and insights from noncommutative Fourier analysis to formulate our algorithms.

The main contributions of our paper include methods for performing general probabilistic inference operations completely in the Fourier domain. In particular, we developed the Kronecker conditioning algorithm, which conditions a distribution on evidence using Bayes rule while operating only on Fourier coefficients. While prediction/rollup operations can be written as pointwise products in the Fourier domain, we showed that conditioning operations can be written, in dual fashion, as generalized convolutions in the Fourier domain. Our conditioning algorithm is general in two senses: first, one can use Kronecker conditioning to handle any observation model which can be written in the Fourier domain, and second, the same algorithm can be applied to condition distributions over arbitrary finite groups. Due to this generality, we are able to efficiently compute the Fourier transforms of a wide variety of probabilistic models which may potentially be useful in different applications.

We presented an analysis of the errors which can accumulate in bandlimited inference and argued that Fourier based approaches work well when the underlying distributions are diffuse and are thus well approximated by low-frequency basis functions. During inference, errors in high-order terms due to bandlimiting can be propagated to lower-order terms and bandlimited conditioning can, on occasion, result in Fourier coefficients which correspond to no valid distribution. We showed, however, that the problem can be remedied by projecting to a relaxation of the marginal polytope.

Finally, our evaluation on data from a camera network shows that our methods perform well when compared to the optimal solution in small problems, or to an omniscient tracker in large problems. Furthermore, we demonstrated that our projection step is fundamental in obtaining these high-quality results.

Algebraic methods have recently enjoyed a surge of interest in the machine learning community. We believe that our unified approach for performing probabilistic inference over permutations, as well as our gentle exposition of group representation theory and noncommutative Fourier analysis will significantly lower the barrier of entry for machine learning researchers who are interested in using or further developing algebraically inspired algorithms which are useful for real-world problems.

## Acknowledgments

This work is supported in part by the Office of Naval Research under MURI N000140710747, the Army Research Office under grant W911NF-06-1-0275, the National Science Foundation under grants DGE-0333420, EEE-540865, NeTS-NOSS 0626151 and TF 0634803, and by the Pennsylvania Infrastructure Technology Alliance (PITA). Carlos Guestrin was also supported in part by an Alfred P. Sloan Fellowship. We are grateful to Kyle Heath for collecting the camera data and to

Robert Hough, Emre Oto and Risi Kondor for helpful and insightful discussions. We would like to thank the reviewers for providing thoughtful and detailed feedback.

## Appendix A. Groups

This section is intended as a quick glossary for the group theoretic definitions used in the paper. Groups are a generalization of many of the spaces that we typically work with, such as the real numbers, integers, vector spaces, and matrices. The definition of a group unifies all of these spaces under a handful of axioms.

**Definition 17 (Group)** *A group is a set  $G$  together with a binary operation  $\cdot : G \times G \rightarrow G$  (called the group operation) such that the following group axioms hold:*

1. (Associativity) *The group operation is associative. That is, for any group elements  $g_1, g_2, g_3 \in G$ , we have:*

$$(g_1 \cdot g_2) \cdot g_3 = g_1 \cdot (g_2 \cdot g_3), \text{ for all } g_1, g_2, g_3 \in G.$$

2. (Identity) *There exists an identity element (denoted by  $\varepsilon$ ) such that  $g \cdot \varepsilon = \varepsilon \cdot g = g$  for any  $g \in G$ .*
3. (Inverses) *For every  $g \in G$ , there exists an inverse element  $g^{-1}$  such that  $g \cdot g^{-1} = g^{-1} \cdot g = \varepsilon$ .*

**Definition 18 (Abelian Group)** *If, for any group elements  $g_1, g_2 \in G$ , we have  $g_1 \cdot g_2 = g_2 \cdot g_1$ , then  $G$  is called an Abelian or commutative group.*

Perhaps the most familiar group is the set of integers,  $\mathbb{Z}$ , with respect to the addition operation. It is well known that for any integers  $a, b, c \in \mathbb{Z}$ ,  $a + (b + c) = (a + b) + c$ . The identity element in the integers is zero, and every element has an additive inverse ( $a + (-a) = (-a) + a = 0$ ). Additionally, the integers are an Abelian group since  $a + b = b + a$  for any  $a, b \in \mathbb{Z}$ . Note that the natural numbers  $\mathbb{N} = \{0, 1, 2, 3, \dots\}$  do not form a group with respect to addition because inverses do not exist.

The main example of a group in this paper, of course, is the symmetric group, the set of permutations of  $\{1, \dots, n\}$ . The group operation on permutations is function composition, which is associative, and we discussed inverses and the identity element in Section 3.

**Example 13** *There are many groups besides the integers and the symmetric group. The following are several examples.*

- *The positive real numbers  $\mathbb{R}^+$  form a group with respect to multiplication. The identity element of  $\mathbb{R}^+$  is the multiplicative identity, 1, and given a real number  $x$ , there exists an inverse element  $\frac{1}{x}$ .*
- *As an example of a finite group, the integers modulo  $n$ ,  $\mathbb{Z}/n\mathbb{Z}$ , form a group with respect to addition modulo  $n$ .*
- *The invertible  $n \times n$  matrices over the reals,  $GL_n(\mathbb{R})$ , form a group with respect to matrix multiplication. The  $n \times n$  identity matrix serves as the identity element in  $GL_n(\mathbb{R})$ , and by assumption, every matrix in  $GL_n(\mathbb{R})$  is invertible.*

The group axioms impose strong structural constraints on  $G$ , and one of the ways that structure is manifested in groups is in the existence of *subgroups*.

**Definition 19 (Subgroup)** If  $G$  is a group (with group operation  $\cdot$ ), a subset  $H \subset G$  is called a subgroup if it is itself a group with respect to the same group operation.  $H$  is called a trivial subgroup if it is either all of  $G$  or consists only of a single element.

**Example 14** We have the following examples of subgroups.

- The even integers,  $2\mathbb{Z}$ , form a subgroup of the integers since the sum of any two even integers is an even integer, and the inverse (negative) of an even integer is again even. However, the odd integers do not form a subgroup since the sum of two odd integers is not odd.
- The special orthogonal matrices (orthogonal matrices with determinant  $+1$ ) form a subgroup of the group of  $n \times n$  matrices,  $GL_n(\mathbb{R})$ . This can be seen by using the facts (1), that  $(\det A)(\det B) = \det(AB)$  and (2), that the inverse of any orthogonal matrix is also orthogonal.

### Appendix B. Constructing Irreducible Representation Matrices

In this section, we present (without proof) some standard algorithms for constructing the irreducible representation matrices with respect to the *Gel'fand-Tsetlin (GZ) basis* (constructed with respect to the subgroup chain  $S_1 \subset S_2 \subset \dots \subset S_n$ ).<sup>13</sup> None of the techniques in Appendix B are novel. For a more elaborate discussion, see, for example, Kondor (2006), Chen (1989) and Vershik and Okounkov (2006). There are several properties which make the irreducible representation matrices, written with respect to the GZ basis, fairly useful in practice. They are guaranteed to be, for example, real-valued and orthogonal. And as we will show, the matrices have certain useful sparsity properties that can be exploited in implementation.

We begin by introducing a few concepts relating to *Young tableaux* which are like Young tabloids with the distinction that the rows are considered as *ordered tuples* rather than *unordered sets*. For example, the following two diagrams are distinct as Young tableaux, but not as Young tabloids:

$$\begin{array}{|c|c|c|} \hline 1 & 2 & 3 \\ \hline 4 & 5 & \\ \hline \end{array} \neq \begin{array}{|c|c|c|} \hline 1 & 3 & 2 \\ \hline 5 & 4 & \\ \hline \end{array} \quad (\text{as Young tableaux}).$$

A Young Tableau  $t$  is said to be *standard* if its entries are increasing to the right along rows and down columns. For example, the set of all standard Young Tableaux of shape  $\lambda = (3, 2)$  is:

$$\left\{ \begin{array}{|c|c|c|} \hline 1 & 3 & 5 \\ \hline 2 & 4 & \\ \hline \end{array}, \begin{array}{|c|c|c|} \hline 1 & 2 & 5 \\ \hline 3 & 4 & \\ \hline \end{array}, \begin{array}{|c|c|c|} \hline 1 & 3 & 4 \\ \hline 2 & 5 & \\ \hline \end{array}, \begin{array}{|c|c|c|} \hline 1 & 2 & 4 \\ \hline 3 & 5 & \\ \hline \end{array}, \begin{array}{|c|c|c|} \hline 1 & 2 & 3 \\ \hline 4 & 5 & \\ \hline \end{array} \right\}. \tag{31}$$

Given a permutation  $\sigma \in S_n$ , one can always apply  $\sigma$  to a Young tableau  $t$  to get a new Young tableau, which we denote by  $\sigma \circ t$ , by permuting the labels within the tableau. For example,

$$(1, 2) \circ \begin{array}{|c|c|c|} \hline 1 & 2 & 3 \\ \hline 4 & 5 & \\ \hline \end{array} = \begin{array}{|c|c|c|} \hline 2 & 1 & 3 \\ \hline 4 & 5 & \\ \hline \end{array}.$$

Note, however, that even if  $t$  is a standard tableau,  $\sigma \circ t$  is not guaranteed to be standard.

The significance of the standard tableaux is that the set of all standard tableaux of shape  $\lambda$  can be used to index the set of GZ basis vectors for the irreducible representation  $\rho_\lambda$ . Since there are

<sup>13</sup> The irreducible representation matrices in this Appendix are also sometimes referred to as *Young's Orthogonal Representation (YOR)*.

five total standard tableaux of shape  $(3, 2)$ , we see, for example, that the irreducible corresponding to the partition  $(3, 2)$  is 5-dimensional. There is a simple recursive procedure for enumerating the set of all standard tableaux of shape  $\lambda$ , which we illustrate for  $\lambda = (3, 2)$ .

**Example 15** If  $\lambda = (3, 2)$ , there are only two possible boxes that the label 5 can occupy so that both rows and columns are increasing. They are:

$$\begin{array}{|c|c|c|} \hline & & 5 \\ \hline & & \\ \hline \end{array}, \text{ and } \begin{array}{|c|c|c|} \hline & & \\ \hline & & 5 \\ \hline \end{array}.$$

To enumerate the set of all standard tableaux of shape  $(3, 2)$ , we need to fill the empty boxes in the above partially filled tableaux with the labels  $\{1, 2, 3, 4\}$  so that both rows and columns are increasing. Enumerating the standard tableaux of shape  $(3, 2)$  thus reduces to enumerating the set of standard tableaux of shapes  $(2, 2)$  and  $(3, 1)$ , respectively. For  $(2, 2)$ , the set of standard tableaux (which, in implementation would be computed recursively) is:

$$\left\{ \begin{array}{|c|c|} \hline 1 & 3 \\ \hline 2 & 4 \\ \hline \end{array}, \begin{array}{|c|c|} \hline 1 & 2 \\ \hline 3 & 4 \\ \hline \end{array} \right\},$$

and for  $(3, 1)$ , the set of standard tableaux is:

$$\left\{ \begin{array}{|c|c|c|} \hline 1 & 3 & 4 \\ \hline 2 & & \\ \hline \end{array}, \begin{array}{|c|c|c|} \hline 1 & 2 & 4 \\ \hline 3 & & \\ \hline \end{array}, \begin{array}{|c|c|c|} \hline 1 & 2 & 3 \\ \hline 4 & & \\ \hline \end{array} \right\}.$$

The entire set of standard tableaux of shape  $(3, 2)$  is therefore:

$$\left\{ \begin{array}{|c|c|c|} \hline 1 & 3 & 5 \\ \hline 2 & 4 & \\ \hline \end{array}, \begin{array}{|c|c|c|} \hline 1 & 2 & 5 \\ \hline 3 & 4 & \\ \hline \end{array} \right\} \cup \left\{ \begin{array}{|c|c|c|} \hline 1 & 3 & 4 \\ \hline 2 & 5 & \\ \hline \end{array}, \begin{array}{|c|c|c|} \hline 1 & 2 & 4 \\ \hline 3 & 5 & \\ \hline \end{array}, \begin{array}{|c|c|c|} \hline 1 & 2 & 3 \\ \hline 4 & 5 & \\ \hline \end{array} \right\}.$$

Before explicitly constructing the representation matrices, we must define a signed distance on Young Tableaux called the *axial distance*.

**Definition 20** The axial distance,  $d_t(i, j)$ , between entries  $i$  and  $j$  in tableau  $t$ , is defined to be:

$$d_t(i, j) \equiv (col(t, j) - col(t, i)) - (row(t, j) - row(t, i)),$$

where  $row(t, i)$  denotes the row of label  $i$  in tableau  $t$ , and  $col(t, i)$  denotes the column of label  $i$  in tableau  $t$ .

Intuitively, the axial distance between  $i - 1$  and  $i$  in a standard tableau  $t$  is equal to the (signed) number of steps that are required to travel from  $i - 1$  to  $i$ , if at each step, one is allowed to traverse a single box in the tableau in one of the four cardinal directions. For example, the axial distance from 3 to 4 with respect to tableau:  $t = \begin{array}{|c|c|c|} \hline 1 & 2 & 3 \\ \hline 4 & 5 & \\ \hline \end{array}$  is:

$$\begin{aligned} d_t(3, 4) &= \left( col \left( \begin{array}{|c|c|c|} \hline 1 & 2 & 3 \\ \hline 4 & 5 & \\ \hline \end{array}, 4 \right) - col \left( \begin{array}{|c|c|c|} \hline 1 & 2 & 3 \\ \hline 4 & 5 & \\ \hline \end{array}, 3 \right) \right) - \left( row \left( \begin{array}{|c|c|c|} \hline 1 & 2 & 3 \\ \hline 4 & 5 & \\ \hline \end{array}, 4 \right) - row \left( \begin{array}{|c|c|c|} \hline 1 & 2 & 3 \\ \hline 4 & 5 & \\ \hline \end{array}, 3 \right) \right) \\ &= (1 - 3) - (2 - 1) = -3 \end{aligned}$$

### B.1 Constructing Representation Matrices for Adjacent Transpositions

In the following discussion, we will consider a fixed ordering,  $t_1, \dots, t_{d_\lambda}$ , on the set of standard tableaux of shape  $\lambda$  and refer to both standard tableaux and columns of  $\rho_\lambda(\sigma)$  interchangeably. Thus  $t_1$  refers to first column,  $t_2$  refers to the second column and so on. And we will index elements in  $\rho_\lambda(\sigma)$  using pairs of standard tableau,  $(t_j, t_k)$ .

To explicitly define the representation matrices with respect to the GZ basis, we will first construct the matrices for adjacent transpositions (i.e., permutations of the form  $(i - 1, i)$ ), and then we will construct arbitrary representation matrices by combining the matrices for the adjacent transpositions. The rule for constructing the matrix coefficient  $[\rho_\lambda(i - 1, i)]_{t_j, t_k}$  is as follows.

1. Define the  $(t_j, t_k)$  coefficient of  $\rho_\lambda(i - 1, i)$  to be zero if it is (1), off-diagonal ( $j \neq k$ ) and (2), not of the form  $(t_j, (i - 1, i) \circ t_k)$ .
2. If  $(t_j, t_k)$  is a diagonal element, (i.e., of the form  $(t_j, t_j)$ ), define:

$$[\rho_\lambda(i - 1, i)]_{t_j, t_j} = 1/d_{t_j}(i - 1, i),$$

where  $d_{t_j}(i - 1, i)$  is the axial distance which we defined earlier in the section.

3. If  $(t_j, t_k)$  can be written as  $(t_j, (i - 1, i) \circ t_j)$  define:

$$[\rho_\lambda(i - 1, i)]_{t_j, \sigma \circ t_j} = \sqrt{1 - 1/d_{t_j}^2(i - 1, i)}.$$

Note that the only time that off-diagonal elements can be nonzero under the above rules is when  $(i - i, i) \circ t_j$  happens to also be a standard tableau. If we apply an adjacent transposition,  $\sigma = (i - 1, i)$  to a standard tableau  $t$ , then  $\sigma \circ t$  is guaranteed to be standard if and only if  $i - 1$  and  $i$  were neither in the same row nor column of  $t$ . This can be seen by examining each case separately.

1.  **$i - 1$  and  $i$  are in the same row or same column of  $t$ .** If  $i$  and  $i - 1$  are in the same row of  $t$ , then  $i - 1$  lies to the left of  $i$ . Applying  $\sigma \circ t$  swaps their positions so that  $i$  lies to the left of  $i - 1$ , and so we see that  $\sigma \circ t$  cannot be standard. For example,

$$(3, 4) \circ \begin{array}{|c|c|c|} \hline 1 & 2 & 5 \\ \hline 3 & 4 & \\ \hline \end{array} = \begin{array}{|c|c|c|} \hline 1 & 2 & 5 \\ \hline 4 & 3 & \\ \hline \end{array}.$$

Similarly, we see that if  $i$  and  $i - 1$  are in the same column of  $t$ ,  $\sigma \circ t$  cannot be standard. For example,

$$(3, 4) \circ \begin{array}{|c|c|c|} \hline 1 & 3 & 5 \\ \hline 2 & 4 & \\ \hline \end{array} = \begin{array}{|c|c|c|} \hline 1 & 4 & 5 \\ \hline 2 & 3 & \\ \hline \end{array}.$$

2.  **$i - 1$  and  $i$  are neither in the same row nor column of  $t$ .** In the second case,  $\sigma \circ t$  can be seen to be a standard tableau due to the fact that  $i - 1$  and  $i$  are adjacent indices. For example,

$$(3, 4) \circ \begin{array}{|c|c|c|} \hline 1 & 2 & 3 \\ \hline 4 & 5 & \\ \hline \end{array} = \begin{array}{|c|c|c|} \hline 1 & 2 & 4 \\ \hline 3 & 5 & \\ \hline \end{array}.$$

Therefore, to see if  $(i - 1, i) \circ t$  is standard, we need only check to see that  $i - 1$  and  $i$  are in different rows and columns of the tableau  $t$ . The pseudocode for constructing the irreducible representation matrices for adjacent swaps is summarized in Algorithm 3. Note that the matrices constructed in the algorithm are sparse, with no more than two nonzero elements in any given column.



**Algorithm 3:** Pseudocode for computing irreducible representations matrices with respect to the Gel'fand-Tsetlin basis at adjacent transpositions.

```

ADJACENTRHO
input :  $i \in \{2, \dots, n\}, \lambda$ 
output:  $\rho_\lambda(i-1, i)$ 
1  $\rho \leftarrow \mathbf{0}_{d_\lambda \times d_\lambda}$ ;
2 foreach standard tableaux  $t$  of shape  $\lambda$  do
3    $d \leftarrow (\text{col}(t, i) - \text{col}(t, i-1)) - (\text{row}(t, i) - \text{row}(t, i-1))$ ;
4    $\rho(t, t) \leftarrow 1/d$ ;
5   if  $i-1$  and  $i$  are in different rows and columns of  $t$  then
6      $\rho((i-1, i) \circ (t), t) \leftarrow \sqrt{1 - 1/d^2}$ ;
7 return  $\rho$ ;
    
```

**Example 16** We compute the representation matrix of  $\rho_{(3,2)}$  evaluated at the adjacent transposition  $\sigma = (i-1, i) = (3, 4)$ . For this example, we will use the enumeration of the standard tableaux of shape  $(3, 2)$  given in Equation 31.

For each  $(3, 2)$ -tableau  $t_j$ , we identify whether  $\sigma \circ t_j$  is standard and compute the axial distance from 3 to 4 on the tableau  $t_j$ .

$j$	1	2	3	4	5
$t_j$	$\begin{array}{ c c c } \hline 1 & 3 & 5 \\ \hline 2 & 4 & \\ \hline \end{array}$	$\begin{array}{ c c c } \hline 1 & 2 & 5 \\ \hline 3 & 4 & \\ \hline \end{array}$	$\begin{array}{ c c c } \hline 1 & 3 & 4 \\ \hline 2 & 5 & \\ \hline \end{array}$	$\begin{array}{ c c c } \hline 1 & 2 & 4 \\ \hline 3 & 5 & \\ \hline \end{array}$	$\begin{array}{ c c c } \hline 1 & 2 & 3 \\ \hline 4 & 5 & \\ \hline \end{array}$
$(3, 4) \circ t_j$	$\begin{array}{ c c c } \hline 1 & 4 & 5 \\ \hline 2 & 3 & \\ \hline \end{array}$	$\begin{array}{ c c c } \hline 1 & 2 & 5 \\ \hline 4 & 3 & \\ \hline \end{array}$	$\begin{array}{ c c c } \hline 1 & 4 & 3 \\ \hline 2 & 5 & \\ \hline \end{array}$	$\begin{array}{ c c c } \hline 1 & 2 & 3 \\ \hline 4 & 5 & \\ \hline \end{array}$	$\begin{array}{ c c c } \hline 1 & 2 & 4 \\ \hline 3 & 5 & \\ \hline \end{array}$
$(3, 4) \circ t_j$ Standard?	No	No	No	Yes	Yes
axial distance ( $d_{t_j}(3, 4)$ )	-1	1	1	3	-3

Putting the results together in a matrix yields:

$$\rho_{(3,2)}(3, 4) = \begin{bmatrix} & t_1 & t_2 & t_3 & t_4 & t_5 \\ t_1 & -1 & & & & \\ t_2 & & 1 & & & \\ t_3 & & & 1 & & \\ t_4 & & & & \frac{1}{3} & \sqrt{\frac{8}{9}} \\ t_5 & & & & \sqrt{\frac{8}{9}} & -\frac{1}{3} \end{bmatrix},$$

where all of the empty entries are zero.

### B.2 Constructing Representation Matrices for General Permutations

To construct representation matrices for general permutations, it is enough to observe that all permutations can be factored into a sequence of adjacent swaps. For example, the permutation  $(1, 2, 5)$  can be factored into:

$$(1, 2, 5) = (4, 5)(3, 4)(1, 2)(2, 3)(3, 4)(4, 5),$$

---

**Algorithm 4:** Pseudocode for computing irreducible representation matrices for arbitrary permutations.

---

```

GETRHO
input :  $\sigma \in S_n, \lambda$ 
output:  $\rho_\lambda(\sigma)$  (a  $d_\lambda \times d_\lambda$  matrix)
1 //Use Bubblesort to factor  $\sigma$  into a product of transpositions
2  $k \leftarrow 0$ ;
3  $factors \leftarrow \emptyset$ ;
4 for  $i = 1, 2, \dots, n$  do
5   for  $j = n, n-1, \dots, i+1$  do
6     if  $\sigma(j) < \sigma(j-1)$  then
7       Swap( $\sigma(j-1), \sigma(j)$ );
8        $k \leftarrow k+1$ ;
9        $factors(k) \leftarrow j$ ;
10 //Construct representation matrix using adjacent transpositions
11  $\rho_\lambda(\sigma) \leftarrow I_{d_\lambda \times d_\lambda}$ ;
12  $m \leftarrow \text{length}(factors)$ ;
13 for  $j = 1, \dots, m$  do
14    $\rho_\lambda(\sigma) \leftarrow \text{GETADJACENTRHO}(factors(j), \lambda) \cdot \rho_\lambda(\sigma)$ ;

```

---

and hence, for any partition  $\lambda$ ,

$$\rho_\lambda(1, 2, 5) = \rho_\lambda(4, 5) \cdot \rho_\lambda(3, 4) \cdot \rho_\lambda(1, 2) \cdot \rho_\lambda(2, 3) \cdot \rho_\lambda(3, 4) \cdot \rho_\lambda(4, 5),$$

since  $\rho_\lambda$  is a group representation. Algorithmically, factoring a permutation into adjacent swaps looks very similar to the Bubblesort algorithm, and we show the pseudocode in Algorithm 4.

### Appendix C. Fourier Transforming the Indicator Function $\delta_{S_{X,Y}}$

In this section, we derive the Fourier transform of the indicator function of the two-sided coset  $S_{X,Y} \subset S_n$  (see Equation 27). To do so, we will need to understand the Gel'fand-Tsetlin basis at a slightly deeper level. For example, the fact that the basis elements are indexed by standard tableaux has not been motivated and may seem unintuitive. We begin this section by motivating the standard tableaux from the perspective of the *branching rule*, a standard fact taken from the representation theory of the symmetric group. We then show how the branching rule leads to a method for computing the desired indicator functions.

#### C.1 Standard Tableaux and the Gel'fand-Tsetlin Basis

It is straightforward to see that any irreducible representation  $\rho_\lambda$  of  $S_n$ , can also be seen as a representation of  $S_{n-1}$  when restricted to permutations in  $S_{n-1}$  (the set of elements which fix  $n$ ). We will denote the restricted representation by  $\rho_\lambda \downarrow_{S_{n-1}}^{S_n}$ . However, the irreducibility property may not be preserved by restriction, which is to say that, as a representation of  $S_{n-1}$ ,  $\rho_\lambda$  is not necessarily irreducible and might decompose as Equation 2 would dictate. Thus, for all  $\sigma \in S_{n-1}$  (or more

precisely,  $\sigma \in S_n$  such that  $\sigma(n) = n$ , there exists  $C_\lambda$  and multiplicities  $z_{\lambda\mu}$  such that:

$$C_\lambda^{-1} \cdot \rho_\lambda(\sigma) \cdot C_\lambda = \bigoplus_{\mu} \bigoplus_{j=1}^{z_{\lambda\mu}} \rho_\mu(\sigma),$$

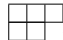
where  $\mu$  ranges over the partitions of  $n - 1$ .

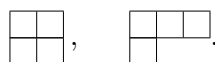
The *branching rule* allows us to state the decomposition even more precisely. Given any partition  $\lambda$  of  $n$ , let  $\lambda^-$  index over the set of partitions of  $n - 1$  whose Ferrers diagrams differ from  $\lambda$  in a single box.

**Theorem 21 (Branching Rule, see Vershik and Okounkov (2006) for a proof)** *For each irreducible representation  $\rho_\lambda$  of  $S_n$ , there exists a matrix  $C_\lambda$  such that:*

$$C_\lambda^{-1} \cdot \rho_\lambda(\sigma) \cdot C_\lambda = \bigoplus_{\lambda^-} \rho_{\lambda^-}(\sigma)$$

holds for any  $\sigma \in S_{n-1}$ .

**Example 17** *If  $\lambda = (3, 2)$ , then its corresponding Ferrers diagram is: , and the Ferrers diagrams corresponding to partitions of 4 which differ from  $\lambda$  in a single box are:*



Thus,  $\lambda^-$  indexes over the set  $\{(2, 2), (3, 1)\}$ . The branching rule states that given an irreducible matrix representation  $\rho_{(3,2)}$  of  $S_5$ , then there is a matrix  $C_{(3,2)}$  such that, for any permutation  $\sigma \in S_5$  such that  $\sigma(5) = 5$ ,

$$C_\lambda^{-1} \cdot \rho_{(3,2)}(\sigma) \cdot C_\lambda = \left[ \begin{array}{c|c} \rho_{(2,2)}(\sigma) & 0 \\ \hline 0 & \rho_{(3,1)}(\sigma) \end{array} \right].$$

The Gel'fand-Tsetlin basis is constructed such that the branching rule holds with all  $C_\lambda = I$ . Thus the irreducible representation matrices constructed with respect to the GZ basis have the property that the equation:

$$\rho_\lambda(\sigma) = \bigoplus_{\lambda^-} \rho_{\lambda^-}(\sigma)$$

holds identically for all  $\sigma \in S_{n-1}$ . We now can show how the branching rule naturally leads to indexing the basis elements by standard tableaux. First observe that the branching rule allows us to associate each column of the irreducible  $\rho_\lambda$  with some partition of  $n - 1$ .

If we recursively apply the branching rule again (thus restricting to  $S_{n-2}$ ), we see that the following decomposition holds:

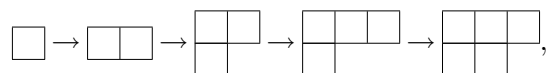
$$\rho_\lambda(\sigma) = \bigoplus_{\lambda^-} \left[ \bigoplus_{\lambda^{--}} \rho_{\lambda^{--}}(\sigma) \right],$$

where  $\lambda^{--}$  indexes over partitions which differ from  $\lambda^-$  by a single box. Thus each column can be associated with a partition of  $n - 1$  and a partition of  $n - 2$ . Taking this logic even further, we can restrict to  $S_{n-3}$ ,  $S_{n-4}$ , and so on until we can restrict no further, associating each column with a sequence of partitions<sup>14</sup>  $\mu_1 \vdash 1, \mu_2 \vdash 2 \dots, \mu_n \vdash n$ , where each partition  $\mu_i$  can be obtained by adding

14. Here we use the  $\lambda \vdash n$  notation to denote the relation that  $\lambda$  is a partition of  $n$ . For example,  $(3, 2, 1) \vdash 6$ .

a single box to the Ferrers diagram of  $\mu_{i-1}$ , and  $\mu_n = \lambda$ . We will refer to such a sequence as a *branching sequence*. Since the branching rule guarantees multiplicity-free decompositions (that is,  $z_{\lambda\mu} = 1$  for all pairs  $(\lambda, \mu)$ ), it turns out that each column of  $\rho_\lambda$  is *uniquely* specified by a branching sequence.

**Example 18** A possible branching sequence is:



or written as partitions,  $[(1) \rightarrow (2) \rightarrow (2, 1) \rightarrow (3, 1) \rightarrow (3, 2)]$ .

The set of all possible branching sequences ending in  $\lambda$  can be visualized using a *branching tree* (shown for  $\lambda = (3, 2)$  in Figure 12(a)), where each branching sequence is a path between the root and some leaf node. We will denote the branching tree corresponding to the partition  $\lambda$  by  $\mathcal{T}^\lambda$  and the set of nodes at the  $r^{\text{th}}$  level of  $\mathcal{T}^\lambda$  by  $\mathcal{T}_r^\lambda$  (where the root node forms the zeroth level by convention). We can rephrase the branching rule in terms of the branching tree.

**Proposition 22** Let  $\rho_\lambda$  be an irreducible matrix representation of  $S_n$  (constructed with respect to the Gel'fand-Tsetlin basis). For any  $\sigma \in S_k \subset S_n$ ,  $\rho_\lambda(\sigma)$  decomposes as:

$$\rho_\lambda(\sigma) = \bigoplus_{\mu \in \mathcal{T}_{n-k}^\lambda} \rho_\mu(\sigma).$$

**Example 19** As an example, consider applying Proposition 22 to  $\rho_{(3,2)}$  with  $k = 3$ . The  $(n - k)$ th (second) level of the branching tree for  $\lambda = (3, 2)$ ,  $\mathcal{T}_2^{(3,2)}$  consists of two copies of the partition  $(2, 1)$  and a single copy of the partition  $(3)$ . Thus for any element  $\sigma \in S_5$  which fixes 4 and 5 ( $\sigma(4) = 4$ ,  $\sigma(5) = 5$ ), we have:

$$\rho_{(3,2)}(\sigma) = \left[ \begin{array}{c|c|c} \rho_{(2,1)}(\sigma) & & \\ \hline & \rho_{(2,1)}(\sigma) & \\ \hline & & \rho_{(3)}(\sigma) \end{array} \right].$$

As a final remark, observe that branching sequences can be compactly represented as *standard tableaux*, where the number in each box indicates the point in the sequence at which the box was added. For example, the following standard tableau and sequence of partitions are equivalent:

$$\begin{array}{|c|c|c|} \hline 1 & 2 & 4 \\ \hline 3 & 5 & \\ \hline \end{array} \longleftrightarrow [(1) \rightarrow (2) \rightarrow (2, 1) \rightarrow (3, 1) \rightarrow (3, 2)].$$

To summarize, the GZ basis (adapted to the subgroup chain  $S_1 \subset \dots \subset S_n$ ) is defined so that the branching rule holds as a matrix identity (with no need of a change of basis matrix), and furthermore, each basis vector of the representation space for  $\rho_\lambda$  can be associated with a branching sequence, or equivalently, a standard tableau.

### C.2 Fourier Transforming $\delta_{S_{X,Y}}$

We are now in a position to compute the Fourier transform of indicators of the form  $\delta_{S_{X,Y}}$ . First, as a corollary of the branching rule, we see that we can decompose the Fourier transform of functions that are supported on  $S_{n-1} \subset S_n$ .

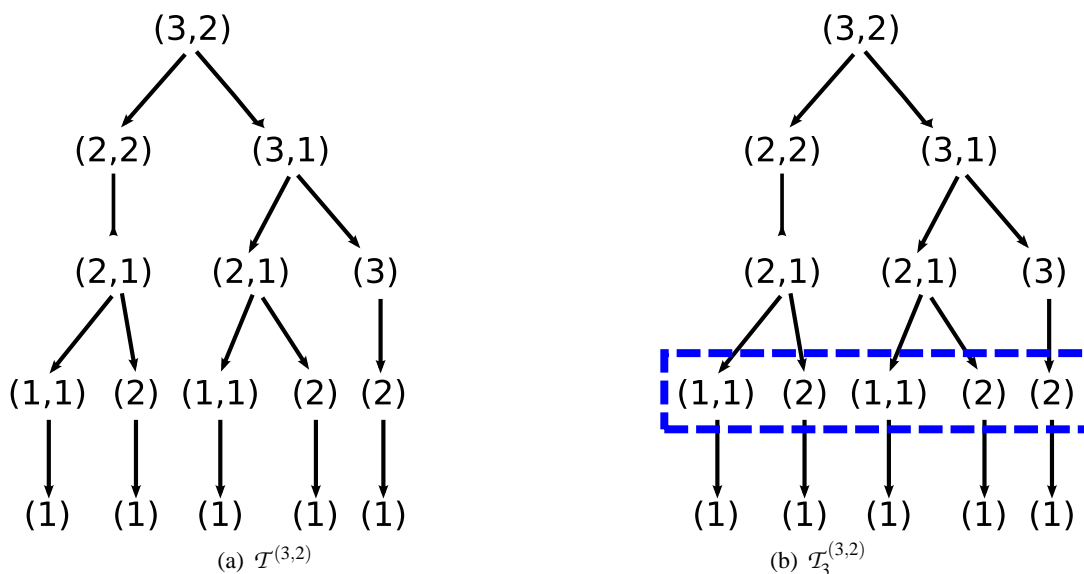


Figure 12: (a) The branching tree for  $\lambda = (3, 2)$ . (b) The 3<sup>rd</sup> level of  $\mathcal{T}^{(3,2)}$  (outlined) is denoted by  $\mathcal{T}_3^{(3,2)}$  and consists of two copies of the partition  $(1, 1)$  and three copies of the partition  $(2)$ .

**Corollary 23** *If  $f : S_n \rightarrow \mathbb{R}$  is supported on the subgroup  $S_{n-1}$ , then for each partition  $\lambda$ , the Fourier transform of  $f$  (with respect to the Gel'fand-Tsetlin basis adapted  $S_1 \subset S_2 \subset \dots \subset S_n$ ) decomposes into a direct sum of Fourier transforms on  $S_{n-1}$ . Specifically, we have:*

$$\hat{f}_\lambda = \bigoplus_{\lambda^-} [\hat{f} \downarrow_{n-1}^n]_{\lambda^-},$$

where  $f \downarrow_{n-1}^n$  is defined to be the restriction of  $f$  to  $S_{n-1}$ .

Consider the Fourier transform of the indicator function of  $S_k \subset S_n$ :

$$\delta_{S_k}(\sigma) = \begin{cases} 1 & \text{if } \sigma(j) = j \text{ for } j \in \{k+1, \dots, n\} \\ 0 & \text{otherwise} \end{cases}.$$

We now apply the branching rule  $n - k$  times to the indicator function  $\delta_{S_k}$ . Since  $\delta_{S_k}$  is supported on  $S_k \subset S_n$ , the Fourier transform of  $\delta_{S_k}$  at the irreducible  $\rho_\lambda$  can be written as a direct sum of Fourier coefficient matrices at the irreducibles which appear in the  $n - k$ th level of the branching tree corresponding to  $\lambda$ .

$$\left[ \hat{\delta}_{S_k} \right]_\lambda = \bigoplus_{\mu \in \mathcal{T}_{n-k}^\lambda} \left[ \hat{\delta}_{S_k} \downarrow_k^n \right]_\mu$$

Furthermore, since the restriction of  $\delta_{S_k}$  to the subgroup  $S_k$  is a constant function, we see that all of the nontrivial irreducible summands are zero (since the Fourier transform of a constant function is zero at all nontrivial terms) and that the trivial terms are exactly  $k!$ . Because the trivial representation is one-dimensional, only a subset of the diagonal elements of  $\left[ \hat{\delta}_{S_k} \right]_\lambda$  can be nonzero.

---

**Algorithm 5:** Pseudocode for computing the Fourier transform of the indicator function of  $S_k \subset S_n$  at the partition  $\lambda$ .

---

$S_k$ -INDICATOR  
**input** :  $k, n, \lambda$  (a partition of  $n$ )  
**output**:  $\left[ \widehat{\delta}_{S_k} \right]_{\lambda}$

- 1  $\left[ \widehat{\delta}_{S_k} \right]_{\lambda} \leftarrow \mathbf{0}_{d_{\lambda} \times d_{\lambda}}$  ;
- 2 **foreach** standard tableaux  $t$  of shape  $\lambda$  **do**
- 3     **if**  $t \downarrow_k^n = \boxed{1} \boxed{2} \boxed{3} \cdots \boxed{k}$  **then**
- 4          $\left[ \widehat{\delta}_{S_k} \right]_{\lambda}(t, t) \leftarrow k!$ ;

---

Algorithmically we can construct the Fourier transform of  $\delta_{S_k}$  at  $\lambda$  by enumerating all of the branching sequences for  $\lambda$  and setting the  $(j, j)$  diagonal element of  $\left[ \widehat{\delta}_{S_k} \right]_{\lambda}$  to be  $k!$  if the corresponding  $j$ th branching sequence contains the partition  $(k)$ . Alternatively, we can state the procedure in terms of standard tableaux. First, we define a restriction operation on a standard tableau  $t$ .

**Definition 24** Given a standard tableau  $t$  with  $n$  boxes and a positive integer  $k < n$ , we define the restriction of  $t$  to  $S_k$  (denoted by  $t \downarrow_k^n$ ) to be the standard tableau  $t$  after removing boxes containing labels  $k + 1, \dots, n$ .

To construct the Fourier transform of  $\delta_{S_k}$  at  $\lambda$ , we iterate through the standard tableaux of shape  $\lambda$ , and set the  $(j, j)$  diagonal element of  $\left[ \widehat{\delta}_{S_k} \right]_{\lambda}$  to be  $k!$  if the restriction of the  $j$ th tableau to  $S_k$ ,  $t_j \downarrow_k^n$ , takes the form  $\boxed{1} \boxed{2} \boxed{3} \cdots \boxed{k}$ . See Algorithm 5.

**Example 20** We compute  $\left[ \widehat{\delta}_{S_2} \right]_{(3,2)}$  as an example. The branching sequences for  $\lambda = (3, 2)$  are:

$$\begin{array}{lcl}
 \begin{array}{|c|c|c|} \hline 1 & 3 & 5 \\ \hline 2 & 4 & \\ \hline \end{array} & \longleftrightarrow & [(1) \rightarrow (1, 1) \rightarrow (2, 1) \rightarrow (2, 2) \rightarrow (3, 2)], \\
 \begin{array}{|c|c|c|} \hline 1 & 2 & 5 \\ \hline 3 & 4 & \\ \hline \end{array} & \longleftrightarrow & [(1) \rightarrow (2) \rightarrow (2, 1) \rightarrow (2, 2) \rightarrow (3, 2)], \\
 \begin{array}{|c|c|c|} \hline 1 & 3 & 4 \\ \hline 2 & 5 & \\ \hline \end{array} & \longleftrightarrow & [(1) \rightarrow (1, 1) \rightarrow (2, 1) \rightarrow (3, 1) \rightarrow (3, 2)], \\
 \begin{array}{|c|c|c|} \hline 1 & 2 & 4 \\ \hline 3 & 5 & \\ \hline \end{array} & \longleftrightarrow & [(1) \rightarrow (2) \rightarrow (2, 1) \rightarrow (3, 1) \rightarrow (3, 2)], \\
 \begin{array}{|c|c|c|} \hline 1 & 2 & 3 \\ \hline 4 & 5 & \\ \hline \end{array} & \longleftrightarrow & [(1) \rightarrow (2) \rightarrow (3) \rightarrow (3, 1) \rightarrow (3, 2)].
 \end{array}$$

Since there are only three sequences which contain the partition  $(2)$ , only those three basis elements have nonzero entries. And finally, noting that the appropriate normalization constant here

is simply  $|S_2| = 2! = 2$ , we see that:

$$\left[ \hat{\delta}_{S_2} \right]_{(3,2)} = \begin{array}{c|ccccc} & \begin{array}{|c|c|c|} \hline 1 & 3 & 5 \\ \hline 2 & 4 & \\ \hline \end{array} & \begin{array}{|c|c|c|} \hline 1 & 2 & 5 \\ \hline 3 & 4 & \\ \hline \end{array} & \begin{array}{|c|c|c|} \hline 1 & 3 & 4 \\ \hline 2 & 5 & \\ \hline \end{array} & \begin{array}{|c|c|c|} \hline 1 & 2 & 4 \\ \hline 3 & 5 & \\ \hline \end{array} & \begin{array}{|c|c|c|} \hline 1 & 2 & 3 \\ \hline 4 & 5 & \\ \hline \end{array} \\ \hline \begin{array}{|c|c|c|} \hline 1 & 3 & 5 \\ \hline 2 & 4 & \\ \hline \end{array} & 0 & 0 & 0 & 0 & 0 \\ \begin{array}{|c|c|c|} \hline 1 & 2 & 5 \\ \hline 3 & 4 & \\ \hline \end{array} & 0 & 2 & 0 & 0 & 0 \\ \begin{array}{|c|c|c|} \hline 1 & 3 & 4 \\ \hline 2 & 5 & \\ \hline \end{array} & 0 & 0 & 0 & 0 & 0 \\ \begin{array}{|c|c|c|} \hline 1 & 2 & 4 \\ \hline 3 & 5 & \\ \hline \end{array} & 0 & 0 & 0 & 2 & 0 \\ \begin{array}{|c|c|c|} \hline 1 & 2 & 3 \\ \hline 4 & 5 & \\ \hline \end{array} & 0 & 0 & 0 & 0 & 2 \end{array}.$$

Our discussion has been focused on the indicator function  $\delta_{S_k}$ , but computing  $\delta_{S_{X,Y}}$  with  $|X| = |Y| = n - k$  can be accomplished by first constructing the Fourier coefficient matrices for  $\delta_{S_k}$ , then relabeling the tracks and identities using a change of basis. More precisely, suppose that, to achieve this relabeling, we must permute the  $X$  (identities) using a permutation  $\pi_1$  and the  $Y$  (tracks) using  $\pi_2$ . The *Shift Theorem* (see Diaconis 1988) can be applied to reorder the Fourier coefficients according to these new labels.

**Proposition 25 (Shift Theorem)** *Given  $f : S_n \rightarrow \mathbb{R}$ , define  $f' : S_n \rightarrow \mathbb{R}$  by  $f'(\sigma) = f(\pi_1 \sigma \pi_2)$  for some fixed  $\pi_1, \pi_2 \in S_n$ . The Fourier transforms of  $f$  and  $f'$  are related as:  $\hat{f}'_\lambda = \rho_\lambda(\pi_1) \cdot \hat{f}_\lambda \cdot \rho_\lambda(\pi_2)$ .*

We conclude with a comment on sparsity. It is clear from Algorithm 5 that the coefficient matrices of  $\left[ \hat{\delta}_{S_{X,Y}} \right]_\lambda$  are all, up to an appropriate relabeling of identities and tracks, diagonal matrices with at most  $O(d_\lambda)$  nonzero entries. In fact, we can sometimes show that a given model has  $O(1)$  nonzero entries.

Consider, for example, the indicator function  $\delta_{S_{n-1}}$ , corresponding to observations of the form (“Identity  $j$  is at track  $i$ ”), which is nonzero only at the first two partitions,  $(n)$ , and  $(n-1, 1)$ . The zeroth-order term is,  $\left[ \hat{\delta}_{S_{n-1}} \right]_{(n)} = (n-1)!$ . The first-order Fourier coefficient matrix,  $\left[ \hat{\delta}_{S_{n-1}} \right]_{(n-1,1)}$ , is a matrix of all zeroes except for a single element on the diagonal,  $\left[ \hat{\delta}_{S_{n-1}} \right]_{(n-1,1)}(t, t)$ , where  $t = \begin{array}{|c|c|c|c|} \hline 1 & 2 & 3 & \dots \\ \hline n & & & \end{array}$ , which takes on the value  $(n-1)!$ .

## Appendix D. Decomposing the Tensor Product Representation

We now turn to the *Tensor Product Decomposition* problem, which is that of finding the irreducible components of the typically reducible tensor product representation. If  $\rho_\lambda$  and  $\rho_\mu$  are irreducible representations of  $S_n$ , then there exists an intertwining operator  $C_{\lambda\mu}$  such that:

$$C_{\lambda\mu}^{-1} \cdot (\rho_\lambda \otimes \rho_\mu(\sigma)) \cdot C_{\lambda\mu} = \bigoplus_{\nu} \bigoplus_{\ell=1}^{z_{\lambda\mu\nu}} \rho_\nu(\sigma). \quad (32)$$

In this section, we will present a set of numerical methods for computing the Clebsch-Gordan series ( $z_{\lambda\mu}$ ) and Clebsch-Gordan coefficients ( $C_{\lambda\mu}$ ) for a pair of irreducible representations  $\rho_\lambda \otimes \rho_\mu$ . We begin by discussing two methods for computing the Clebsch-Gordan series. In the second section, we provide a general algorithm for computing the intertwining operators which relate two equivalent representations and discuss how it can be applied to computing the Clebsch-Gordan coefficients (Equation 32) and the matrices which relate marginal probabilities to irreducible Fourier coefficients (Equation 6). The results of Appendix D.1 are specific to the symmetric group, while the results of Appendix D.2 can be applied to arbitrary finite groups.

### D.1 Computing the Clebsch-Gordan Series

We begin with a simple, well-known algorithm based on *group characters* for computing the Clebsch-Gordan series that turns out to be computationally intractable, but yields several illuminating theoretical results. See Serre (1977) for proofs of the theoretical results cited in this section.

One of the main results of representation theory was the discovery that there exists a relatively compact way of encoding any representation up to equivalence with a vector which we call the *character* of the representation. If  $\rho$  is a representation of a group  $G$ , then the character of the representation  $\rho$ , is defined simply to be the trace of the representation at each element  $\sigma \in G$ :

$$\chi_\rho(\sigma) = \text{Tr}(\rho(\sigma)).$$

The reason characters have been so extensively studied is that they uniquely characterize a representation up to equivalence in the sense that two characters  $\chi_{\rho_1}$  and  $\chi_{\rho_2}$  are equal if and only if  $\rho_1$  and  $\rho_2$  are equivalent as representations. Even more surprising is that the space of possible group characters is orthogonally spanned by the characters of the irreducible representations. To make this precise, we first define an inner product on functions from  $G$ .

**Definition 26** Let  $\phi, \psi$  be two real-valued functions on  $G$ . The inner product of  $\phi$  and  $\psi$  is defined to be:

$$\langle \phi, \psi \rangle \equiv \frac{1}{|G|} \sum_{\sigma \in G} \phi(\sigma)\psi(\sigma)$$

With respect to the above inner product, we have the following important result which allows us to test a given representation for irreducibility, and to test two irreducibles for equivalence.

**Proposition 27** Let  $\chi_{\rho_1}$  and  $\chi_{\rho_2}$  be characters corresponding to irreducible representations. Then

$$\langle \chi_{\rho_1}, \chi_{\rho_2} \rangle = \begin{cases} 1 & \text{if } \rho_1 \equiv \rho_2 \\ 0 & \text{otherwise} \end{cases}.$$

Proposition 27 shows that the irreducible characters form an orthonormal set of functions. The next proposition says that the irreducible characters *span* the space of all possible characters.

**Proposition 28** Suppose  $\rho$  is any representation of  $G$  and which decomposes into irreducibles as:

$$\rho \equiv \bigoplus_{\lambda} \bigoplus_{\ell=1}^{z_\lambda} \rho_\lambda,$$

where  $\lambda$  indexes over all irreducibles of  $G$ . Then:



1. The character of  $\rho$  is a linear combination of irreducible characters ( $\chi_\rho = \sum_\lambda z_\lambda \chi_{\rho_\lambda}$ ),
2. and the multiplicity of each irreducible,  $z_\lambda$ , can be recovered using  $\langle \chi_\rho, \chi_{\rho_\lambda} \rangle = z_\lambda$ .

A simple way to decompose any group representation  $\rho$ , is given by Proposition 28, which says that we can take inner products of  $\chi_\rho$  against the basis of irreducible characters to obtain the irreducible multiplicities  $z_\lambda$ . To treat the special case of finding the Clebsch-Gordan series, one observes that the character of the tensor product is simply the pointwise product of the characters of each tensor product factor.

**Theorem 29** *Let  $\rho_\lambda$  and  $\rho_\mu$  be irreducible representations with characters  $\chi_\lambda, \chi_\mu$  respectively. Let  $z_{\lambda\mu\nu}$  be the number of copies of  $\rho_\nu$  in  $\rho_\lambda \otimes \rho_\mu$  (hence, one term of the Clebsch-Gordan series). Then:*

1. The character of the tensor product representation is given by:

$$\chi_{\rho_\lambda \otimes \rho_\mu} = \chi_\lambda \cdot \chi_\mu = \sum_{\nu} z_{\lambda\mu\nu} \chi_\nu.$$

2. The terms of the Clebsch-Gordan series can be computed using:

$$z_{\lambda\mu\nu} = \frac{1}{|G|} \sum_{g \in G} \chi_\lambda(g) \cdot \chi_\mu(g) \cdot \chi_\nu(g),$$

and satisfy the following symmetry:

$$z_{\lambda\mu\nu} = z_{\lambda\nu\mu} = z_{\mu\lambda\nu} = z_{\mu\nu\lambda} = z_{\nu\lambda\mu} = z_{\nu\mu\lambda}. \quad (33)$$

Dot products for characters on the symmetric group can be done in  $O(\#(n))$  time where  $\#(n)$  is the number of partitions of the number  $n$ , instead of the naive  $O(n!)$  time. In practice however,  $\#(n)$  also grows too quickly for the character method to be tractable.

#### D.1.1 MURNAGHAN'S FORMULAS

A theorem by Murnaghan (1938) gives us a ‘bound’ on which representations can appear in the tensor product decomposition on  $S_n$ .

**Theorem 30** *Let  $\rho_1, \rho_2$  be the irreducibles corresponding to the partition  $(n - p, \lambda_2, \dots)$  and  $(n - q, \mu_2, \dots)$  respectively. Then the product  $\rho_1 \otimes \rho_2$  does not contain any irreducibles corresponding to a partition whose first term is less than  $n - p - q$ .*

In view of the connection between the Clebsch-Gordan series and convolution of Fourier coefficients, Theorem 30 is analogous to the fact that for functions over the reals, the convolution of two compactly supported functions is also compactly supported.

We can use Theorem 30 to show that Kronecker conditioning is exact at certain irreducibles.

**Proof** [of Theorem 12] Let  $\Lambda$  denote the set of irreducibles at which our algorithm maintains Fourier coefficients. Since the errors in the prior come from setting coefficients outside of  $\Lambda$  to be zero, we see that Kronecker conditioning returns an approximate posterior which is exact at the irreducibles in

$$\Lambda_{EXACT} = \{\rho_\nu : z_{\lambda\mu\nu} = 0, \text{ where } \lambda \notin \Lambda \text{ and } \mu \succeq (n - q, \mu_2, \dots)\}.$$

Combining Theorem 30 with Equation 33: if  $z_{\lambda\mu\nu} > 0$ , with  $\lambda = (n - p, \lambda_2, \lambda_3, \dots)$ ,  $\mu = (n - q, \mu_2, \mu_3, \dots)$  and  $\nu = (n - r, \nu_2, \nu_3, \dots)$ , then we have that:  $r \leq p + q$ ,  $p \leq q + r$ , and  $q \leq p + r$ . In particular, it implies that  $r \geq p - q$  and  $r \geq q - p$ , or more succinctly,  $r \geq |p - q|$ . Hence, if  $\nu = (n - r, \nu_2, \dots)$ , then  $\rho_\nu \in \Lambda_{EXACT}$  whenever  $r \leq |p - q|$ , which proves the desired result. ■

The same paper (Murnaghan, 1938) derives several general Clebsch-Gordan series formulas for pairs of low-order irreducibles in terms of  $n$ , and in particular, derives the Clebsch-Gordan series for many of the Kronecker product pairs that one would likely encounter in practice. For example,

- $\rho_{(n-1,1)} \otimes \rho_{(n-1,1)} \equiv \rho_{(n)} \oplus \rho_{(n-1,1)} \oplus \rho_{(n-2,2)} \oplus \rho_{(n-2,1,1)}$
- $\rho_{(n-1,1)} \otimes \rho_{(n-2,2)} \equiv \rho_{(n-1,1)} \oplus \rho_{(n-2,2)} \oplus \rho_{(n-2,1,1)} \oplus \rho_{(n-3,3)} \oplus \rho_{(n-3,2,1)}$
- $\rho_{(n-1,1)} \otimes \rho_{(n-2,1,1)} \equiv \rho_{(n-1,1)} \oplus \rho_{(n-2,2)} \oplus \rho_{(n-2,1,1)} \oplus \rho_{(n-3,2,1)} \oplus \rho_{(n-3,1,1,1)}$
- $\rho_{(n-1,1)} \otimes \rho_{(n-3,3)} \equiv \rho_{(n-2,2)} \oplus \rho_{(n-3,3)} \oplus \rho_{(n-3,2,1)} \oplus \rho_{(n-4,4)} \oplus \rho_{(n-4,3,1)}$

### D.2 Computing the Clebsch-Gordan Coefficients

In this section, we consider the general problem of finding an orthogonal operator which decomposes an arbitrary complex representation,  $X(\sigma)$ , of a finite group  $G$ .<sup>15</sup> Unlike the Clebsch-Gordan series which are basis-independent, intertwining operators must be recomputed if we change the underlying basis by which the irreducible representation matrices are constructed. However, for a fixed basis, we remind the reader that these intertwining operators need only be computed once and for all and can be stored in a table for future reference. Let  $X$  be any degree  $d$  group representation of  $G$ , and let  $Y$  be an equivalent direct sum of irreducibles, for example,

$$Y(\sigma) = \bigoplus_{\nu} \bigoplus_{\ell=1}^{z_{\nu}} \rho_{\nu}(\sigma), \tag{34}$$

where each irreducible  $\rho_{\nu}$  has degree  $d_{\nu}$ . We would like to compute an invertible (and orthogonal) operator  $C$ , such that  $C \cdot X(\sigma) = Y(\sigma) \cdot C$ , for all  $\sigma \in G$ . Throughout this section, we will assume that the multiplicities  $z_{\nu}$  are known. To compute Clebsch-Gordan coefficients, for example, we would set  $X = \rho_{\lambda} \otimes \rho_{\mu}$ , and the multiplicities would be given by the Clebsch-Gordan series (Equation 32). To find the matrix which relates marginal probabilities to irreducible coefficients, we would set  $X = \tau_{\lambda}$ , and the multiplicities would be given by the Kostka numbers (Equation 6).

We will begin by describing an algorithm for computing a basis for the space of all possible intertwining operators which we denote by:

$$\text{Int}_{[X;Y]} = \{C \in \mathbb{R}^{d \times d} : C \cdot X(\sigma) = Y(\sigma) \cdot C, \forall \sigma \in G\}.$$

We will then discuss some of the theoretical properties of  $\text{Int}_{[X;Y]}$  and show how to efficiently select an *orthogonal* element of  $\text{Int}_{[X;Y]}$ .

---

15. Though the fundamental ideas in this section hold for a general finite group, we will continue to index irreducible by partitions and think of representations as being real-valued. To generalize the results, one can simply replace all transposes in this section by adjoints and think of  $\eta$  as indexing over the irreducibles of  $G$  rather than partitions.

Our approach is to naively<sup>16</sup> view the task of finding elements of  $\text{Int}_{[X;Y]}$  as a similarity matrix recovery problem, with the twist that the similarity matrix must be consistent over all group elements. To the best of our knowledge, the technique presented in this section is original. We first cast the problem of recovering a similarity matrix as a nullspace computation.

**Proposition 31** *Let  $A, B, C$  be matrices and let  $K_{AB} = I \otimes A - B^T \otimes I$ . Then  $AC = CB$  if and only if  $\text{vec}(C) \in \text{Nullspace}(K_{AB})$ .*

**Proof** A well known matrix identity (van Loan, 2000) states that if  $A, B, C$  are matrices, then  $\text{vec}(ABC) = (C^T \otimes A) \text{vec}(B)$ . Applying the identity to  $AC = CB$ , we have:

$$\text{vec}(ACI) = \text{vec}(ICB),$$

and after some manipulation:

$$(I \otimes A - B^T \otimes I) \text{vec}(C) = 0,$$

showing that  $\text{vec}(C) \in \text{Nullspace}(K_{AB})$ . ■

For each  $\sigma \in G$ , the nullspace of the matrix  $K(\sigma)$  constructed using the above proposition as:

$$K(\sigma) = I \otimes Y(\sigma) - X(\sigma) \otimes I, \quad (35)$$

where  $I$  is a  $d \times d$  identity matrix, corresponds to the space of matrices  $C_\sigma$  such that

$$C_\sigma \cdot X(\sigma) = Y(\sigma) \cdot C, \quad \text{for all } \sigma \in G.$$

To find the space of intertwining operators which are consistent across all group elements, we need to find the intersection:

$$\bigcap_{\sigma \in G} \text{Nullspace}(K(\sigma)).$$

At first glance, it may seem that computing the intersection might require examining  $n!$  nullspaces if  $G = S_n$ , but as luck would have it, most of the nullspaces in the intersection are extraneous, as we now show.

**Definition 32** *We say that a finite group  $G$  is generated by a set of generators  $S = \{g_1, \dots, g_m\}$  if every element of  $G$  can be written as a finite product of elements in  $S$ .*

For example, the following three sets are all generators for  $S_n$ :

- $\{(1,2), (1,3), \dots, (1,n)\}$ ,
- $\{(1,2), (2,3), (3,4), \dots, (n-1,n)\}$ , and
- $\{(1,2), (1,2,3, \dots, n)\}$ .

To ensure a consistent similarity matrix for all group elements, we use the following proposition which says that it suffices to be consistent on any set of generators of the group.

---

16. In implementation, we use a more efficient algorithm for computing intertwining operators known as the *Eigenfunction Method* (EFM) (Chen, 1989). Unfortunately, the EFM is too complicated for us to describe in this paper. The method which we describe in this appendix is conceptually simpler than the EFM and generalizes easily to groups besides  $S_n$ .

**Proposition 33** *Let  $X$  and  $Y$  be representations of finite group  $G$  and suppose that  $G$  is generated by the elements  $\sigma_1, \dots, \sigma_m$ . If there exists an invertible linear operator  $C$  such that  $C \cdot X(\sigma_i) = Y(\sigma_i) \cdot C$  for each  $i \in \{1, \dots, m\}$ , then  $X$  and  $Y$  are equivalent as representations with  $C$  as the intertwining operator.*

**Proof** We just need to show that  $C$  is a similarity transform for any other element of  $G$  as well. Let  $\pi$  be any element of  $G$  and suppose  $\pi$  can be written as the following product of generators:  $\pi = \prod_{i=1}^n \sigma_i$ . It follows that:

$$\begin{aligned} C^{-1} \cdot Y(\pi) \cdot C &= C^{-1} \cdot Y\left(\prod_i \sigma_i\right) \cdot C = C^{-1} \cdot \left(\prod_i Y(\sigma_i)\right) \cdot C \\ &= (C^{-1} \cdot Y(\sigma_1) \cdot C)(C^{-1} \cdot Y(\sigma_2) \cdot C) \cdots (C^{-1} \cdot Y(\sigma_m) \cdot C) \\ &= \prod_i (C^{-1} \cdot Y(\sigma_i) \cdot C) = \prod_i X(\sigma_i) = X\left(\prod_i \sigma_i\right) = X(\pi). \end{aligned}$$

Since this holds for every  $\pi \in G$ , we have shown  $C$  to be an intertwining operator between the representations  $X$  and  $Y$ . ■

The good news is that despite having  $n!$  elements,  $S_n$  can be generated by just two elements, namely,  $(1, 2)$  and  $(1, 2, \dots, n)$ , and so the problem reduces to solving for the intersection of two nullspaces,  $(K(1, 2) \cap K(1, 2, \dots, n))$ , which can be done using standard numerical methods. Typically, the nullspace is multidimensional, showing that, for example, the Clebsch-Gordan coefficients for  $\rho_\lambda \otimes \rho_\mu$  are not unique even up to scale.

Because  $\text{Int}_{[X;Y]}$  contains singular operators (the zero matrix is a member of  $\text{Int}_{[X;Y]}$ , for example), not every element of  $\text{Int}_{[X;Y]}$  is actually a legitimate intertwining operator as we require invertibility. In practice, however, since the singular elements correspond to a measure zero subset of  $\text{Int}_{[X;Y]}$ , one method for reliably selecting an operator from  $\text{Int}_{[X;Y]}$  that “works” is to simply select a random element from the nullspace to be  $C$ . It may, however, be desirable to have an *orthogonal* matrix  $C$  which works as an intertwining operator. In the following, we discuss an object called the *Commutant Algebra* which will lead to several insights about the space  $\text{Int}_{[X;Y]}$ , and in particular, will lead to an algorithm for ‘modifying’ any invertible intertwining operator  $C$  to be an *orthogonal* matrix.

**Definition 34** *The Commutant Algebra of a representation  $Y$  is defined to be the space of operators which commute with  $Y$ :<sup>17</sup>*

$$\text{Com}_Y = \{S \in \mathbb{R}^{d \times d} : S \cdot Y(\sigma) = Y(\sigma) \cdot S, \forall \sigma \in G\}.$$

The elements of the Commutant Algebra of  $Y$  can be shown to always take on a particular constrained form (shown using Schur’s Lemma in Sagan 2001). In particular, every element of  $\text{Com}_Y$  takes the form

$$S = \bigoplus_v (M_{z_v} \otimes I_{d_v}), \tag{36}$$

where  $M_{z_v}$  is some  $z_v \times z_v$  matrix of coefficients and  $I_{d_v}$  is the  $d_v \times d_v$  identity (recall that the  $z_v$  are the multiplicities from Equation 34). Moreover, it can be shown that every matrix of this form must necessarily be an element of the Commutant Algebra.

<sup>17</sup>. Notice that the definition of the Commutant Algebra does not involve the representation  $X$ .

The link between  $\text{Com}_Y$  and our problem is that the space of intertwining operators can be thought of as a ‘translate’ of the Commutant Algebra.

**Lemma 35** *There exists a vector space isomorphism between  $\text{Int}_{[X;Y]}$  and  $\text{Com}_Y$ .*

**Proof** Let  $R$  be any invertible element of  $\text{Int}_{[X;Y]}$  and define the linear map  $f : \text{Com}_Y \rightarrow \mathbb{R}^{d \times d}$  by:  $f : S \mapsto (S \cdot R)$ . We will show that the image of  $f$  is exactly the space of intertwining operators. Consider any element  $\sigma \in G$ :

$$\begin{aligned} (S \cdot R) \cdot X(\sigma) \cdot (S \cdot R)^{-1} &= S \cdot R \cdot X(\sigma) \cdot R^{-1} \cdot S^{-1}, \\ &= S \cdot Y(\sigma) \cdot S^{-1} \quad (\text{since } R \in \text{Int}_{[X;Y]}), \\ &= Y(\sigma) \quad (\text{since } S \in \text{Com}_Y). \end{aligned}$$

We have shown that  $S \cdot R \in \text{Int}_{[X;Y]}$ , and since  $f$  is linear and invertible, we have that  $\text{Int}_{[X;Y]}$  and  $\text{Com}_Y$  are isomorphic as vector spaces. ■

Using the lemma, we can see that the dimension of  $\text{Int}_{[X;Y]}$  must be the same as the dimension of  $\text{Com}_Y$ , and therefore we have the following expression for the dimension of  $\text{Int}_{[X;Y]}$ .

**Proposition 36**

$$\dim \text{Int}_{[X;Y]} = \sum_{\nu} z_{\nu}^2.$$

**Proof** To compute the dimension of  $\text{Int}_{[X;Y]}$ , we need to compute the dimension of  $\text{Com}_Y$ , which can be accomplished simply by computing the number of free parameters in Equation 36. Each matrix  $M_{z_{\nu}}$  is free and yields  $z_{\nu}^2$  parameters, and summing across all irreducibles  $\nu$  yields the desired dimension. ■

To select an orthogonal intertwining operator, we will assume that we are given some invertible  $R \in \text{Int}_{[X;Y]}$  which is not necessarily orthogonal (such as a random element of the nullspace of  $K$ , Equation 35). To find an orthogonal element, we will ‘modify’  $R$  to be an orthogonal matrix by applying an appropriate rotation, such that  $R \cdot R^T = I$ . We begin with a simple observation about  $R \cdot R^T$ .

**Lemma 37** *If both  $X$  and  $Y$  are orthogonal representations and  $R$  is an invertible member of  $\text{Int}_{[X;Y]}$ , then the matrix  $R \cdot R^T$  is an element of  $\text{Com}_Y$ .*

**Proof** Consider a fixed  $\sigma \in G$ . Since  $R \in \text{Int}_{[X;Y]}$ , we have that:

$$X(\sigma) = R^{-1} \cdot Y(\sigma) \cdot R.$$

It is also true that:

$$X(\sigma^{-1}) = R^{-1} \cdot Y(\sigma^{-1}) \cdot R. \quad (37)$$

Since  $X(\sigma)$  and  $Y(\sigma)$  are orthogonal matrices by assumption, Equation 37 becomes:

$$X^T(\sigma) = R^{-1} \cdot Y^T(\sigma) \cdot R.$$

---

**Algorithm 6:** Pseudocode for computing an orthogonal intertwining operators

---

INTXY

**input** : A degree  $d$  orthogonal matrix representation  $X$  evaluated at permutations  $(1, 2)$  and  $(1, \dots, n)$ , and the multiplicity  $z_v$ , of the irreducible  $\rho_v$  in  $X$

**output:** A matrix  $C_v$  with orthogonal rows such that  $C_v^T \cdot \bigoplus^{z_v} \rho_v \cdot C_v = X$

- 1  $K_1 \leftarrow I_{d \times d} \otimes (\bigoplus^{z_v} \rho_v(1, 2)) - X(1, 2) \otimes I_{d \times d}$ ;
  - 2  $K_2 \leftarrow I_{d \times d} \otimes (\bigoplus^{z_v} \rho_v(1, \dots, n)) - X(1, \dots, n) \otimes I_{d \times d}$ ;
  - 3  $K \leftarrow [K_1; K_2]$ ; //Stack  $K_1$  and  $K_2$
  - 4  $v \leftarrow \text{SparseNullspace}(K, z_v^2)$ ; //Find the  $d_v^2$ -dimensional nullspace
  - 5  $R \leftarrow \text{Reshape}(v; z_v d_v, d)$ ; //Reshape  $v$  into a  $(z_v d_v) \times d$  matrix
  - 6  $M \leftarrow \text{KroneckerFactors}(R \cdot R^T)$ ; //Find  $M$  such that  $R \cdot R^T = M \otimes I_{d_v}$
  - 7  $S_v \leftarrow \text{Eigenvectors}(M)$  ;
  - 8  $C_v \leftarrow S_v^T \cdot R$  ;
  - 9 **NormalizeRows**( $C_v$ );
- 

Taking transposes,

$$X(\sigma) = R^T \cdot Y(\sigma) \cdot (R^{-1})^T.$$

We now multiply both sides on the left by  $R$ , and on the right by  $R^T$ ,

$$\begin{aligned} R \cdot X(\sigma) \cdot R^T &= R \cdot R^T \cdot Y(\sigma) \cdot (R^{-1})^T \cdot R^T \\ &= R \cdot R^T \cdot Y(\sigma). \end{aligned}$$

Since  $R \in \text{Int}_{[X;Y]}$ ,

$$Y(\sigma) \cdot R \cdot R^T = R \cdot R^T \cdot Y(\sigma),$$

which shows that  $R \cdot R^T \in \text{Com}_Y$ . ■

We can now state and prove our orthogonalization procedure, which works by diagonalizing the matrix  $R \cdot R^T$ . Due to its highly constrained form, the procedure is quite efficient.

**Theorem 38** *Let  $X$  be any orthogonal group representation of  $G$  and  $Y$  an equivalent orthogonal irreducible decomposition (As in Equation 34). Then for any invertible element  $R \in \text{Int}_{[X;Y]}$ , there exists an (efficiently computable) orthogonal matrix  $T$  such that the matrix  $T \cdot R$  is an element of  $\text{Int}_{[X;Y]}$  and is orthogonal.*

**Proof** Lemma 37 and Equation 36 together imply that the matrix  $R \cdot R^T$  can always be written in the form

$$R \cdot R^T = \bigoplus_v (M_{z_v} \otimes I_{d_v})$$

Since  $R \cdot R^T$  is symmetric, each of the matrices  $M_{z_v}$  is also symmetric and must therefore possess an orthogonal basis of eigenvectors. Define the matrix  $S_{z_v}$  to be the matrix whose columns are the eigenvectors of  $M_{z_v}$ .

The matrix  $S = \bigoplus_v (S_{z_v} \otimes I_{d_v})$  has the following two properties:

1.  $(S^T \cdot R)(S^T \cdot R)^T$  is a diagonal matrix:

Each column of  $S$  is an eigenvector of  $R \cdot R^T$  by standard properties of the direct sum and Kronecker product. Since each of the matrices,  $S_{z_v}$ , is orthogonal, the matrix  $S$  is also orthogonal. We have:

$$\begin{aligned} (S^T \cdot R)(S^T \cdot R)^T &= S^T \cdot R \cdot R^T \cdot S, \\ &= S^{-1} \cdot R \cdot R^T \cdot S, \\ &= D, \end{aligned}$$

where  $D$  is a diagonal matrix of eigenvalues of  $R \cdot R^T$ .

2.  $S^T \cdot R \in \text{Int}_{[X;Y]}$ :

By Equation 36, a matrix is an element of  $\text{Com}_Y$  if and only if it takes the form  $\oplus_v (S_{z_v} \otimes I_{d_v})$ . Since  $S$  can be written in the required form, so can  $S^T$ . We see that  $S^T \in \text{Com}_Y$ , and by the proof of Lemma 35, we see that  $S^T \cdot R \in \text{Int}_{[X;Y]}$ .

Finally, setting  $T = D^{1/2} \cdot S^T$  makes the matrix  $T \cdot R$  orthogonal (and does not change the fact that  $T \cdot R \in \text{Int}_{[X;Y]}$ ). ■

We see that the complexity of computing  $T$  is dominated by the eigenspace decomposition of  $M_{z_v}$ , which is  $O(z_v^3)$ . Pseudocode for computing orthogonal intertwining operators is given Algorithm 6.

## References

- Hamsa Balakrishnan, Inseok Hwang, and Claire Tomlin. Polynomial approximation algorithms for belief matrix maintenance in identity management. In *Proceedings of the 43rd IEEE Conference on Decision and Control, Bahamas*, 2004.
- Yaakov Bar-Shalom and Thomas E. Fortmann. *Tracking and Data Association*. Academic Press, 1988.
- John Bartholdi, Craig Tovey, and Michael Trick. Voting schemes for which it can be difficult to tell who won. *Social Choice and Welfare*, 6(2), 1989.
- Xavier Boyen and Daphne Koller. Tractable inference for complex stochastic processes. In *UAI '98: Uncertainty in Artificial Intelligence*, 1998.
- Jin-Quan Chen. *Group Representation Theory for Physicists*. World Scientific, 1989.
- Michael Clausen and Ulrich Baum. Fast Fourier transforms for symmetric groups: Theory and implementation. *Mathematics of Computations*, 61(204):833–847, 1993.
- Joseph Collins and Jeffrey Uhlmann. Efficient gating in data association with multivariate distributed states. *IEEE Transactions Aerospace and Electronic Systems*, 28, 1992.
- James Cooley and John Tukey. An algorithm for the machine calculation of complex fourier series. *Mathematical Computation*, 19:297–301, 1965.

- Ingemar Cox and Sunita Hingorani. An efficient implementation of Reid’s multiple hypothesis tracking algorithm and its evaluation for the purpose of visual tracking. In *International Conference on Pattern Recognition*, pages 437–443, 1994.
- Douglas E. Critchlow. *Metric Methods for Analyzing Partially Ranked Data*. Springer-Verlag, 1985.
- Persi Diaconis. *Group Representations in Probability and Statistics*. Institute of Mathematical Statistics, 1988.
- Persi Diaconis. A generalization of spectral analysis with application to ranked data. *The Annals of Statistics*, 17(3):949–979, 1989.
- Michael Fligner and Joseph Verducci. Distance based ranking models. *Journal of the Royal Statistical Society*, 48, 1986.
- Richard Foote, Gagan Mirchandani, and Dan Rockmore. Two-dimensional wreath product transforms. *Journal of Symbolic Computation*, 37(2):187–207, 2004.
- Yoav Freund, Raj Iyer, Robert E. Schapire, and Yoram Singer. An efficient boosting algorithm for combining preferences. *Journal of Machine Learning Research (JMLR)*, 4:933–969, 2003. ISSN 1533-7928.
- Ulf Grenander. *Probabilities on Algebraic Structures*. Wiley, 1963.
- David P. Helmbold and Manfred K. Warmuth. Learning permutations with exponential weights. In *COLT ’07: The Twentieth Annual Conference on Learning Theory*, 2007.
- Jonathan Huang, Carlos Guestrin, and Leonidas Guibas. Efficient inference for distributions on permutations. In *NIPS ’07: Advances in Neural Information Processing Systems*, Vancouver, Canada, December 2007.
- Jonathan Huang, Carlos Guestrin, Xiaoye Jiang, and Leonidas Guibas. Exploiting probabilistic independence for permutations. In *AISTATS ’09: Artificial Intelligence and Statistics*, Clearwater Beach, Florida, April 2009.
- Gordon James and Adelbert Kerber. *The Representation Theory of the Symmetric Group*. Addison-Wesley, 1981.
- Thorsten Joachims. Optimizing search engines using clickthrough data. In *KDD ’02: Proceedings of the eighth ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 133–142, New York, NY, USA, 2002. ACM.
- Risi Kondor.  $\mathbb{S}_n\text{ob}$ : a C++ library for fast Fourier transforms on the symmetric group, 2006. Available at <http://www.cs.columbia.edu/~risi/Snob/>.
- Risi Kondor and Karsten M. Borgwardt. The skew spectrum of graphs. In *ICML ’08: Proceedings of the 25th International Conference on Machine Learning*, pages 496–503, 2008.
- Risi Kondor, Andrew Howard, and Tony Jebara. Multi-object tracking with representations of the symmetric group. In *AISTATS ’07: Artificial Intelligence and Statistics*, 2007.



- Ka-Lam Kueh, Timothy Olson, Dan Rockmore, and Ki-Seng Tan. Nonlinear approximation theory on finite groups. Technical Report PMA-TR99-191, Department of Mathematics, Dartmouth College, 1999.
- Serge Lang. *Algebra*. Addison-Wesley, 1965.
- Guy Lebanon and Yi Mao. Non-parametric modeling of partially ranked data. In John C. Platt, Daphne Koller, Yoram Singer, and Sam Roweis, editors, *NIPS '07: Advances in Neural Information Processing Systems*, pages 857–864, Cambridge, MA, 2008. MIT Press.
- Colin Mallows. Non-null ranking models. *Biometrika*, 44, 1957.
- David Maslen. The efficient computation of Fourier transforms on the symmetric group. *Mathematics of Computation*, 67:1121–1147, 1998.
- Marina Meila, Kapil Phadnis, Arthur Patterson, and Jeff Bilmes. Consensus ranking under the exponential model. Technical Report 515, University of Washington, Statistics Department, April 2007.
- Francis Murnaghan. The analysis of the kronecker product of irreducible representations of the symmetric group. *American Journal of Mathematics*, 60(3):761–784, 1938.
- Katta G. Murty. An algorithm for ranking all the assignments in order of increasing cost. *Operations Research*, 16:682–687, 1968.
- Songhwai Oh and Shankar Sastry. A polynomial-time approximation algorithm for joint probabilistic data association. In *Proceedings of the American Control Conference, Portland, OR, 2005*.
- Songhwai Oh, Stuart Russell, and Shankar Sastry. Markov chain Monte Carlo data association for general multiple-target tracking problems. In *Proceedings of the IEEE International Conference on Decision and Control, Paradise Island, Bahamas, 2004*.
- Aubrey B. Poore. Multidimensional assignment and multitarget tracking. In *Partitioning Data Sets*, volume 19, pages 169–196. DIMACS Series in Discrete Mathematics and Theoretical Computer Science, 1995.
- Lawrence Rabiner. A tutorial on hidden markov models and selected applications in speech recognition. *Proceedings of the IEEE*, 77(2):257–286, 1989.
- Donald Reid. An algorithm for tracking multiple targets. *IEEE Transactions on Automatic Control*, 6:843–854, 1979.
- Daniel N. Rockmore. The FFT: An algorithm the whole family can use. *Computing in Science and Engineering*, 02(1):60–64, 2000.
- Bruce E. Sagan. *The Symmetric Group*. Springer, April 2001. ISBN 0387950672.
- Brad Schumitsch, Sebastian Thrun, Gary Bradski, and Kunle Olukotun. The information-form data association filter. In *NIPS '05: Advances in Neural Information Processing Systems*, Cambridge, MA, 2005. MIT Press.

- Brad Schumitsch, Sebastian Thrun, Leonidas Guibas, and Kunle Olukotun. The identity management Kalman filter (imkf). In *RSS '06: Proceedings of Robotics: Science and Systems*, Philadelphia, PA, USA, August 2006.
- Jean-Pierre Serre. *Linear Representations of Finite Groups*. Springer-Verlag, 1977.
- Jaewon Shin, Leonidas Guibas, and Feng Zhao. A distributed algorithm for managing multi-target identities in wireless ad-hoc sensor networks. In *IPSN '03: Information Processing in Sensor Networks*, 2003.
- Jaewon Shin, Nelson Lee, Sebastian Thrun, and Leonidas Guibas. Lazy inference on object identities in wireless sensor networks. In *IPSN '05: Information Processing in Sensor Networks*, 2005.
- Michael Taylor, John Guiver, Stephen Robertson, and Tom Minka. Sofrank: optimizing non-smooth rank metrics. In *WSDM '08: Proceedings of the international conference on Web search and web data mining*, pages 77–86, New York, NY, USA, 2008. ACM.
- Audrey Terras. *Fourier Analysis on Finite Groups and Applications*. London Mathematical Society, 1999.
- Jack van Lint and Richard M. Wilson. *A Course in Combinatorics*. Cambridge University Press, 2001.
- Charles F. van Loan. The ubiquitous kronecker product. *Journal of Computational and Applied Mathematics*, 123(1-2):85–100, 2000. ISSN 0377-0427.
- Anatoly Vershik and Andrei Okounkov. A new approach to the representation theory of symmetric groups. ii. *Journal of Mathematical Sciences*, 131(2):5471–5494, 2006.
- Alan Willsky. On the algebraic structure of certain partially observable finite-state markov processes. *Information and Control*, 38:179–212, 1978.