



# Learning Preferences with Side Information

 Vivek F. Farias,<sup>a</sup> Andrew A. Li<sup>b</sup>
<sup>a</sup>Sloan School of Management, Massachusetts Institute of Technology, Cambridge, Massachusetts 02142; <sup>b</sup>Operations Research Center, Massachusetts Institute of Technology, Cambridge, Massachusetts 02142

 Contact: [vivekf@mit.edu](mailto:vivekf@mit.edu),  <http://orcid.org/0000-0002-5856-9246> (VFF); [aali@mit.edu](mailto:aali@mit.edu),  <http://orcid.org/0000-0002-9552-6421> (AAL)

Received: October 8, 2016

Revised: May 1, 2017; July 19, 2017

Accepted: August 14, 2017

 Published Online in Articles in Advance:  
 May 8, 2019

<https://doi.org/10.1287/mnsc.2018.3092>

Copyright: © 2019 INFORMS

**Abstract.** Product and content personalization is now ubiquitous in e-commerce. There are typically not enough available transactional data for this task. As such, companies today seek to use a variety of information on the interactions between a product and a customer to drive personalization decisions. We formalize this problem as one of recovering a large-scale matrix with side information in the form of additional matrices of conforming dimension. Viewing the matrix we seek to recover and the side information we have as slices of a tensor, we consider the problem of *slice recovery*, which is to recover specific slices of “simple” tensors from noisy observations of the entire tensor. We propose a definition of simplicity that on the one hand elegantly generalizes a standard generative model for our motivating problem and on the other hand subsumes low-rank tensors for a variety of existing definitions of tensor rank. We provide an efficient algorithm for slice recovery that is practical for massive data sets and provides a significant performance improvement over state-of-the-art incumbent approaches to tensor recovery. Furthermore, we establish near-optimal recovery guarantees that, in an important regime, represent an order improvement over the best available results for this problem. Experiments on data from a music streaming service demonstrate the performance and scalability of our algorithm.

**History:** Accepted by Noah Gans, stochastic models and simulation.

**Supplemental Material:** The e-companion is available at <https://doi.org/10.1287/mnsc.2018.3092>.

**Keywords:** [personalization](#) • [e-commerce](#) • [online retail](#) • [recommender systems](#) • [collaborative filtering](#) • [matrix recovery](#) • [tensor recovery](#) • [side information](#) • [multi-interaction data](#)

## 1. Introduction

Consider the problem of learning the propensity of individual customers for different products. This is an important challenge in e-commerce, as customer preferences are a core input into a number of operational and marketing activities. For example, recommender systems are pervasive throughout e-commerce where they serve as a tool to drive sales and consumption by decreasing search costs for users. These systems rely on estimates of customer–product propensities to automatically suggest products to users that they are likely to enjoy. User–item propensities are also used to offer personalized experiences in more general service settings. Search results and loyalty programs may all be personalized to a user. One-to-one marketing, including banner advertisements and targeted sales promotions, benefits from accurate estimates of user–item propensities. The list goes on.

As an example of the task above, consider the problem faced by a retailer estimating the likelihood a specific customer might purchase a given product using historical transaction data available across all customers and products in the retailer’s offering. This task is already quite challenging because of its massive scale: Amazon.com has upwards of  $10^8$  active users

and products (Grey 2015, Statista 2016). In addition to scale, however, this task represents a *statistical* challenge: the probability of a given customer purchasing a given product is small, so that the observed matrix of transactions has a very small number of nonzero entries. For instance, in the case of Amazon, the average active user makes *less than a single* transaction in a month (Harris 2015). We will see that this is, in fact, the key problem: meaningful recovery of the underlying probabilities from their corresponding realizations is hard when these probabilities are small.

Ultimately, the only remedy available in this situation is acquiring more data, perhaps beyond historical transactions. In the retail setting, this is indeed feasible: e-commerce businesses are able to capture data from a variety of distinct *types* of observations at the customer–product granularity. For example, besides sales transactions, online retailers record users’ browse and search histories (capturing “browse” and “search” interactions, respectively, between user–product pairs), clickstream data (capturing a “click” interaction), and responses to advertisements and promotions (capturing a “promotion” interaction), just to name a few. The ultimate task, of course, remains predicting the likelihood that a customer will purchase

a given product. Although these “slices” of data encode interactions between a customer and a product that are distinct from a transaction, they may well inform the likelihood of a transaction and ultimately help resolve the challenge of limited data.

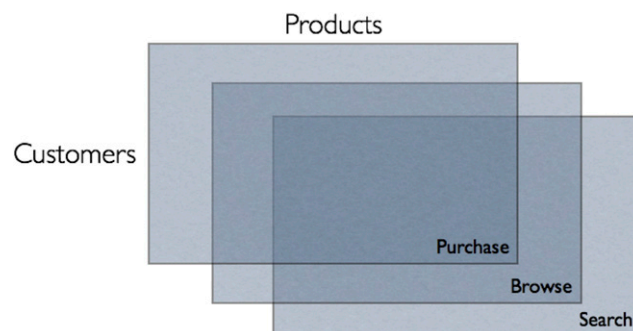
We refer to this as the problem of learning preferences (e.g., the likelihood of a given customer purchasing a given product in the retail example) with side information (e.g., data beyond just historical transactions, as alluded to above). Our first objective will be to formalize this problem; there are many paths we can take here, and we will ultimately propose viewing the problem at hand as one of recovering a three-dimensional *tensor* from its noisy observations. We will then seek to solve this problem and will ultimately present an efficient, near-optimal algorithm for the same that yields a dramatic improvement over existing approaches to tensor recovery.

### 1.1. Representing Data as a Tensor

Taking a step back, the problem of estimating customer–product purchase probabilities from interaction data of a *single* type can be formulated as a *noisy matrix recovery problem*: the goal is to estimate the matrix whose rows correspond to customers, columns correspond to products, and entries contain the purchase probabilities. If our data consisted purely of transactions, these transactions could be viewed as a “noisy realization” of the underlying matrix which we seek to recover. As we discuss in a subsequent section, this formalization can be viewed as equivalent to assigning latent feature vectors  $u_i, v_j$  to each customer  $i$  and each product  $j$ , respectively, and assuming the probability of a specific customer  $i$  purchasing product  $j$  is a bilinear form in these latent vectors,  $f(u_i, v_j)$ . As one allows the dimension of the latent space to grow, the expressive power of such a model grows as well, ultimately becoming fully general. The problem of matrix recovery can be viewed as one of learning the latent feature vectors and the bilinear form  $f(\cdot, \cdot)$  from observed data. Indeed, this formalization of the problem has proved to be incredibly productive and central to the current state of the art in the design of personalization algorithms, as we discuss in the literature review.

Our goal here is to consider multiple types of interactions (e.g., sales transactions *along with* other customer–product interactions, such as browse and search). In analogy to matrix recovery, we may represent these data as a set of matrices, with each matrix corresponding to a single type of interaction. We model the various slices of data that we have as the slices of a three-dimensional tensor (see Figure 1). Tensors, which are the higher-dimensional analogues of matrices, have been used in a broad array of applications to represent high-dimensional data. The underlying generative model (that we will describe in greater detail in the next section) will then, in analogy to

**Figure 1.** (Color online) Customer–Product Interaction Matrices Represented as Slices of a Tensor



the matrix setting, continue to correspond to assigning latent feature vectors  $u_i, v_j$  to each customer  $i$  and each product  $j$ , respectively. However, the likelihood of each interaction type  $k$  between a specific customer  $i$  and product  $j$  will be described by a distinct bilinear form in these latent vectors,  $f_k(u_i, v_j)$ . The problem of “recovering” this tensor from its noisy observations is then equivalent to learning these latent feature vectors along with each of the bilinear forms associated with a specific type of interaction.

We may thus formalize the task we have laid out as the problem of recovering a three-dimensional tensor from its “noisy observations.” In fact, we will typically be interested in recovering a single *slice* of the tensor (for instance, the slice corresponding to the likelihood of a transaction) using data available across *all* slices. To clarify the observation model at hand, the observed data, as a special case, can be seen as a single realization of an underlying tensor of probabilities; that is, a given entry of the observed tensor is a Bernoulli random variable with mean equal to the corresponding entry in the ground truth tensor. When the underlying ground truth probabilities are small, the observed tensor will be sparse (in the usual mathematical sense of the term). This observation model is particularly relevant to the retail example discussed above where the occurrence (or nonoccurrence) of a transaction or other interactions may be thought of as Bernoulli realizations with specific, unknown probabilities. A similar generative model is also relevant to our experiments with music streaming data. This setting, which we will continue to refer to as the problem of *tensor recovery from noisy observations*, or *tensor recovery* for short, is the primary focus of our algorithm and recovery guarantees.

On the other hand, a distinct but closely related observation model is the setting where some small subset of entries of the tensor is observed exactly, and nothing is observed outside of this subset of entries. One application of this setting is the oft-cited Netflix problem: using a data set of user–movie ratings, the task was to predict how users might have rated movies

they had not watched. The data here are naturally represented as a matrix, with some entries observed in the form of ratings and the other entries unobserved. Of course, we can again incorporate side information in this setting, in the form of additional partially observed matrices; this is typically referred to as the *tensor completion* problem. In our retail application, this observation model would correspond to knowing the exact probability of purchase on some subset of customers and nothing outside that subset. Evidently, the observation model is less relevant to the retail application: we do not observe probabilities; we observe transactions. Moreover, we observe whether or not a given customer has purchased a given product for all customers and products; transactions, albeit sparse, are not hidden from the retailer. Nonetheless, we will describe how our algorithm can be extended to this observation model using a device originally proposed by Achlioptas and McSherry (2007).

Before summarizing our approach and contributions, we point out that there is by now a fairly robust literature on the problem of tensor recovery that one might hope to fall back on at this stage. The mainstay of this literature is a convex optimization approach that seeks to find a tensor that is simultaneously “close” to the observed data and “simple” in the sense that a convex surrogate of the “tensor rank” is small. We will formalize the notion of tensor rank in the next section, but ultimately, this quantity can be thought of as restricting the dimension of the latent space of customer and product feature vectors, as well as the family of bilinear functions  $f_k(\cdot, \cdot)$ . Unfortunately, we will see that the convex approach falls short of expectations here. This happens for two key reasons that we will formalize in a subsequent section but explain in brief here. First and foremost, in the tensor setting, these convex optimization approaches are *difficult to scale* to massive amounts of data and will typically call for dense matrix operations at scales that are untenable. This is unlike the matrix setting where the convex approach can be shown equivalent to a simple specialized algorithm (singular value decomposition with soft-thresholding) ideally suited for massive data sets. More important, though, the *statistical power* of these approaches appears to fall well short of what one might hope for in the tensor setting. As we will show in what follows, in our setting of three-dimensional tensors, the error rates achieved using these approaches are akin to what one would get by simply running a matrix recovery algorithm on each individual slice of the tensor (ignoring all other slices!). In addition, *no* guarantees are available on the error of an individual slice, which is particularly relevant because our original motivation is recovering a single slice (for instance, the likelihood of a transaction) using data from all slices.

## 1.2. Our Approach and Contributions

Our approach consists of a simple algorithm for learning the slices of a three-dimensional tensor; applied to all slices, this also results in an algorithm for the recovery of entire tensors from their noisy observations. Relative to the extant literature, we make the following contributions.

### 1.2.1. Statistical Power and Near-Optimal Rates for Noisy Tensor Recovery.

Under a broad set of assumptions, we establish that our approach admits near-optimal guarantees on the rate of recovery for a broad class of three-dimensional tensors from their noisy observations. We accomplish this by establishing upper bounds on the estimation error incurred by using our approach as well as minimax lower bounds applicable to any approach. The guarantees we establish are stronger than those available for existing convex approaches. We simultaneously place looser restrictions on the underlying “ground truth” tensor and the nature of the noise than those required for rigorous recovery via those convex approaches.

Our analysis also admits guarantees on error rates for individual slices that do not have a counterpart in the extant tensor recovery literature. From a pragmatic perspective, this is particularly important in that such guarantees are relevant to our original motivation of recovering a single, specific slice while utilizing data across all slices.

Perhaps the most important aspect of our theoretical guarantees is that they quantify the precise extent to which side information can help with dealing with the problem of sparse data. In fact, a special case of our results establishes a broad set of conditions under which the recovery error on any given single slice (and by implication, the entire tensor) *decays linearly in the number of slices*. Colloquially, this is equivalent to establishing how one may trade off sparsity in data of one type (say, transactions) for data of a different type (say, search data).

### 1.2.2. Empirical Evaluation.

We conduct an extensive set of experiments that empirically demonstrate the most important advantages of our approach. In the first suite of experiments using synthetic data, we benchmark our recovery algorithm against a well-studied algorithm from the family of convex approaches. We observe that the recovery rate of our algorithm exactly matches that predicted by our theoretical results, and similarly, the recovery rate of the convex algorithm matches the best-known theoretical guarantee for convex approaches, which is weaker by an order of magnitude. Moreover, these experiments confirm that our approach is drastically more efficient from a computational standpoint.

In our second suite of experiments, we use real-world data from Xiami.com, a major online music

streaming service. Much like the sales transaction data we have alluded to, these data are extremely sparse, but our algorithm is able to leverage side information to outperform two common benchmarks by a significant margin. These experiments also demonstrate the scalability of our approach in practice. In particular, they are at a large enough scale that most existing tensor recovery algorithms, including convex approaches, are intractable. On the other hand, our algorithm is able to leverage data sparsity to operate on these large-sized tensors.

**1.2.3. Scalability.** In addition to the above, our approach is provably computationally efficient and, more importantly, easy to scale to massive data sets. Specifically, every step within our algorithm can be implemented as a sparse matrix-vector operation and can thus scale to gigantic amounts of data using off-the-shelf computational tools. As a result, we find that our algorithm is typically faster than *a single iteration* of iterative “operator splitting” algorithms used by incumbent approaches, and it is considerably simpler to implement.

We will also show that, with a minor modification, our approach can also apply to the *tensor completion* setting. Here, the nature of the guarantee asks for the number of observed entries needed for eventual recovery as the dimensions of the tensor scale. We will see there that although our guarantee is not order optimal (it is dominated by an alternative but computationally intractable approach), it is still the case that our approach can leverage side information. Specifically, the number of observations needed *per slice* decreases as more slices are added to our approach.

### 1.3. Related Work

Our work falls into various diverse streams of literature relating to the nature of the data used, intended applications, and methodology. We describe these in the following subsections.

**1.3.1. Practical Applications.** The goal of the present paper is to use data to learn user–item propensities, which are a fundamental input to many operational and marketing activities. Perhaps the quintessential application is recommendation. There has been a great amount of work in designing recommender systems; see Ansari et al. (2000) and the references in the nice survey by Adomavicius and Tuzhilin (2005). Recommendations—and, more generally, personalization—are used in all sorts of e-commerce activities and come in many different forms. For some examples in retail, see Linden et al. (2003) for a description of the various ways that Amazon.com recommends products to customers, including targeted emails and shopping cart recommendations.

Outside of retail, Ghose et al. (2012) study hotel recommendations in the form of personalized results in online search engines. Chung et al. (2009) design an adaptive playlist of recommended songs, which is now fairly common in streaming media (e.g., Spotify’s “Discover Weekly” playlist).

The general problem of learning users’ preferences extends to operational activities as well. For example, websites can drive usage, and therefore advertising revenue, through recommendations: Ansari and Mela (2003) study personalized emails, and Besbes et al. (2015) study in-page recommendations to other pages. Fleder and Hosanagar (2009) analyze the impact of recommender systems on the overall demand for all products, and Demirezen and Kumar (2016) look at ways in which recommendations and inventory should be considered jointly. Recommendations have also been studied in the context of promotions (Garfinkel et al. 2008) and the firm’s profit (Hosanagar et al. 2008). Finally, Bodapati (2008) and Jacobs et al. (2016) seek to predict future purchase probabilities; these estimates can have important applications such as demand estimation.

Our work presents a framework for analyzing multiple, diverse data sources, specifically with user–item interaction data. Technological advances have allowed modern businesses to easily record and store massive amounts of data; Naik et al. (2008) survey the many sources of data and point out that the scale and diversity of data collected are constantly expanding. See Section A of the online appendix for a survey of articles dealing with user–item interaction data.

**1.3.2. Statistical Methodology.** In terms of methodology, our work broadly belongs to the class of collaborative filtering algorithms. Since their first introduction by Goldberg et al. (1992), collaborative filtering algorithms have become a mainstay approach in recommender systems. Our work fits within the matrix factorization approaches to these problems (Koren et al. 2009).

A common criticism of collaborative filtering algorithms is their inability to work with sparse data (Ansari et al. 2000). Our approach combines various sources of data to alleviate this sparsity. This fits within a recent stream of research into collaborative filtering algorithms that incorporate more data; see Shi et al. (2014) for a survey of these approaches. One such approach is matrix completion with side information, for example, by Xu et al. (2013), Jain and Dhillon (2013), and Chiang et al. (2015). In this line of work, “side information” refers to user and item features such as user demographics and product specifications. By contrast, our approach centers on interactions between users and items, but we will see that it is flexible enough to incorporate user and item features as well.

Another approach is collective matrix factorization (Singh and Gordon 2008), which studies matrix recovery using multiple matrices across multiple groups. In the setting of multiple matrices across two groups, our approach will, in fact, apply to a more general version of this problem and simultaneously allow us to leverage the modeling advantages of tensors.

Tensors have received significant attention recently in an attempt to generalize matrix recovery to higher dimensions; for a nice survey of tensor decompositions and methods, see Kolda and Bader (2009). Tensors have been applied in a wide range of areas—for example, 3D image and video modeling (Liu et al. 2013), multivariate temporal data analysis (Bahadori et al. 2014), and multitask learning (Romera-Paredes et al. 2013). See Mørup (2011) for a survey of more applications.

We will provide a thorough review of the tensor recovery literature to Section 3.1; for now, we briefly state that the majority of theoretical work so far has been on convex optimization formulations akin to nuclear norm minimization for matrix recovery. This approach was originally proposed by Liu et al. (2013) and Gandy et al. (2011), and it has been analyzed extensively (see Tomioka et al. 2011, Mu et al. 2013, Romera-Paredes and Pontil 2013, and Zheng and Tomioka 2015). Finally, in parallel to the noisy recovery problem, algorithms and guarantees have been shown for the tensor completion problem by Jain and Oh (2014), Huang et al. (2015), and Yuan and Zhang (2016). The goal there is typically to provide conditions under which the challenging goal of *exact* recovery is possible, and those papers prove results under incoherence conditions similar to those made for matrix completion (see Candès and Tao 2010, Gross 2011, and Recht 2011).

Before proceeding, we make the disclaimer that we are focusing solely on three-dimensional tensors, so although our algorithm and guarantee may compare favorably against existing tensor recovery approaches, those approaches are specified for higher-order tensors, whereas ours is not. Alternatively, there is a vast array of applications, including spatiotemporal analysis (Bahadori et al. 2014), neuroimaging (Zhou et al. 2013), and the applications we analyze in the present work, that makes studying 3D tensors specifically important.

## 2. Model and Problem

One common problem that we have already alluded to is predicting customers' preferences for products from sales transaction data. In beginning to formalize this problem, let us label the customers  $i = 1, \dots, m_1$  and products  $j = 1, \dots, m_2$ . To reduce notation, we will set  $m_1 = m_2 = m$ , but the entire analysis in what follows easily generalizes to a different number of customers and products. Say that the data we have are sales transactions that occurred over some previous time period—say, the

previous month. We encode these data in an  $m \times m$  matrix called  $X$ , whose  $(i, j)$ th element  $X_{ij}$  is a binary indicator that equals 1 if customer  $i$  purchased product  $j$  in the last month and 0 otherwise. To formulate a meaningful estimation problem, let us imagine that  $X_{ij}$  is a Bernoulli random variable with mean  $M_{ij}$ , and denote by  $M$  the matrix of these probabilities. More generally, we assume

$$X_{ij} = M_{ij} + \epsilon_{ij},$$

where  $\epsilon_{ij}$  is a mean-zero noise term, and the noise terms are independent of each other but not necessarily identically distributed. To represent Bernoulli observations in this framework, we can let  $\epsilon_{ij} = \text{Ber}(M_{ij}) - M_{ij}$ . Our estimation problem is that of estimating  $M$ , having observed  $X$ .

Making progress on this problem clearly requires structural assumptions on  $M$ ; estimating a completely general such matrix will require a prohibitive amount of data in a sense we will make precise shortly. One fairly general generative model for  $M$  is as follows: let us assume that every customer  $i$  is associated with some unknown vector of *latent* features,  $u_i \in \mathbb{R}^r$ , and every product  $j$  is similarly associated with an unknown latent vector  $v_j \in \mathbb{R}^r$ . We may then consider a generative model of the form

$$M_{i,j} = f(u_i, v_j),$$

where  $f: \mathbb{R}^r \times \mathbb{R}^r \rightarrow \mathbb{R}$  is also unknown. If  $f(\cdot, \cdot)$  were a bilinear form, so that  $f(u_i, v_j) = u_i^\top S v_j$  for some unknown  $S$ , then we may write  $M$  as

$$M = USV^\top,$$

where  $U \in \mathbb{R}^{m \times r}$  has as its  $i$ th row the vector  $u_i^\top$ , and similarly for  $V$ . Our task is now one of estimating a general rank  $r$  matrix  $M$  having observed  $X$ ; the complexity of the underlying generative model is governed by  $r$ . For instance, allowing  $r = m$  permits a fully general model but, as we have noted, one that will be prohibitive to estimate from data. This generative model has been used extensively in operations and marketing to model product purchases (Grover and Srinivasan 1987), ratings (Ansari et al. 2000), and click-throughs (Ansari and Mela 2003), along with applications in myriad settings spanning information retrieval (Berry et al. 1995) and latent semantic analysis (Papadimitriou et al. 1998), computer vision problems such as facial recognition (Sirovich and Kirby 1987) and background subtraction (Oliver et al. 2000), sensor network localization (So and Ye 2007), bioinformatics (Troyanskaya et al. 2001), social network analysis (Liben-Nowell and Kleinberg 2007), and web link analysis (Kleinberg 1999), just to name a few.

The above setup demonstrates the manner in which the problem of predicting customers' preferences for

products using data of a single type (e.g., transactions) can be formulated precisely as one of low-rank matrix recovery, which in turn is an incredibly well-studied problem. It is well known that the minimax mean squared error of any matrix estimator is  $\Omega(r/m)$ . That is, for any estimator  $\hat{M}(\cdot)$ , we can construct a rank  $r$  matrix  $M$  for which

$$\frac{1}{m^2} E\|\hat{M}(X) - M\|_F^2 = \Omega(r/m).$$

In fact, estimators that achieve this lower bound have been constructed (e.g., Koltchinskii et al. 2011). In light of these error rates, recovery is only meaningful in a relative sense if  $\|M\|_F^2$  is larger in size than this error (i.e., if  $\|M\|_F^2 = \Omega(rm)$ ). Now, if  $M$  encodes purchase probabilities (i.e., the entries of  $M$  are in  $[0, 1]$ ), then  $\|M\|_1 \geq \|M\|_F^2$ , so that for recovery to be meaningful, we require  $\|M\|_1 = \Omega(rm)$ . In other words, we *require on average  $r$  expected transactions per user to estimate a rank  $r$  model*. Put another way, the sparsity of observed transactions limits the complexity of the model we can estimate. This limitation is often observed in practice: consider that in 2014, the online music streaming service Spotify had  $10^7$  active users and songs (i.e.,  $m \sim 10^7$ ) and had observed  $10^{10}$  user–song pairs, such that a user listened to a particular song (i.e.,  $\|M\|_1 \sim 10^{10}$ ). The most complex models Spotify estimated for use in its recommendation engine had rank  $\sim 10^3$  (Bernhardsson 2014), which matches the error bound just described. This error bound also makes clear the challenge faced by, say, an online retailer where, as we have noted previously, on average only a single transaction is observed per user. It is this challenge that our work is meant to address.

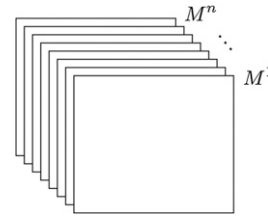
### 2.1. Multiple Data Types and Tensors

Now suppose the data we have come from observing multiple types of customer–product interactions. We label the different types of interactions  $k = 1, \dots, n$ , and we denote the value generated by an interaction of type  $k$  between customer  $i$  and product  $j$  as  $X_{ij}^k$ . As before, these data are a noisy observation of some underlying ground truth value denoted by  $M_{ij}^k$  that is never observed; for example,  $X_{ij}^k$  might be a Bernoulli random variable with mean  $M_{ij}^k$ . In general, we assume

$$X_{ij}^k = M_{ij}^k + \epsilon_{ij}^k,$$

where  $\epsilon_{ij}^k$  is a mean-zero noise term, and the noise terms are independent of each other but not necessarily identically distributed. To model  $X_{ij}^k$  as a Bernoulli random variable with mean  $M_{ij}^k$ , one simply takes  $\epsilon_{ij}^k = \text{Ber}(M_{ij}^k) - M_{ij}^k$ . We let  $M^k$  and  $X^k$  denote the  $m \times m$  matrices whose  $(i, j)$ th elements are  $M_{ij}^k$  and  $X_{ij}^k$ , respectively. Returning to our running example of

**Figure 2.** Matrices  $M^1, \dots, M^n$  Represented as Slices of a Tensor



predicting sales using data on a variety of distinct interactions, in our notation,  $X^1$  may now be the matrix whose  $(i, j)$ th entry is 1 if and only if customer  $i$  purchased product  $j$ ,  $X^2$  could be the matrix whose  $(i, j)$ th entry is 1 if and only if customer  $i$  browsed product  $j$ , and so on. Matrices  $M^1$  and  $M^2$  would then encode the probabilities of those events that we seek to estimate. Further matrices  $X^3, \dots, X^n$  would encode data observed from ratings, search, etc., taken as noisy observations of the corresponding matrices  $M^3, \dots, M^n$ .

It will be convenient to introduce some basic notation for three-dimensional tensors. Let  $M \in \mathbb{R}^{m \times m \times n}$  denote the three-dimensional tensor obtained by stacking the matrices  $M^1, \dots, M^n$  (see Figure 2). We call each of these matrices *slices* of the tensor, and continuing with the same notation, we denote the  $(i, j)$ th element of the  $k$ th slice as  $M_{ij}^k$ . Without loss of generality, we will assume that every entry of  $M$  lies in  $[-1, 1]$ ; this can always be satisfied by rescaling.

A common set of operations on tensors that we will use frequently here is that of *unfolding*, which essentially flattens the tensor into a matrix (see Figure 3). The mode 1 unfolding of  $M$ , denoted by  $M_{(1)}$ , is the  $m \times nm$  matrix whose columns are the columns of  $M^1, \dots, M^n$  (the order of the columns will not matter for us). Similarly, the mode 2 unfolding, denoted by  $M_{(2)}$ , is the  $m \times nm$  matrix whose columns are the transposed rows of  $M^1, \dots, M^n$ .<sup>1</sup>

### 2.2. Generative Model and Slice Rank

Just as in matrix recovery, we cannot hope to recover the tensor  $M$  without some assumption on its simplicity. In analogy with the setting of a single interaction (where we dealt with matrices), we now seek to propose a natural generative model for  $M$ . To that end, we assume that every customer  $i$  is associated with some unknown vector of *latent* features,  $u_i \in \mathbb{R}^r$ , and every product  $j$  is similarly associated with an unknown latent vector,  $v_j \in \mathbb{R}^r$ . We may then consider that for each type of interaction  $k = 1, 2, \dots, n$ , we have

$$M_{ij}^k = f^k(u_i, v_j),$$

where  $f^k: \mathbb{R}^r \times \mathbb{R}^r \rightarrow \mathbb{R}$  is also unknown. This model says that the likelihood of an interaction between a specific user and product depends on feature vectors specific to

**Figure 3.** Graphical Representation of the Mode 1 and Mode 2 Unfoldings of  $M$

$$M_{(1)} = \begin{bmatrix} M^1 & M^2 & \cdots & M^n \end{bmatrix}$$

$$M_{(2)} = \begin{bmatrix} M^{1\top} & M^{2\top} & \cdots & M^{n\top} \end{bmatrix}$$

the user and product, respectively, through a function that is specific to the interaction. Given such a model, the hope is that data from any interaction can now contribute to learning the underlying user and product feature vectors, although the function determining a specific interaction is pinned down using data from just that interaction. It is this intuition that will eventually guide our recovery algorithms.

As before, if  $f^k(\cdot, \cdot)$  were a bilinear form, so that  $f^k(u_i, v_j) = u_i^\top S^k v_j$ , then we may write each of the slices of  $M$  as

$$M^k = US^kV^\top,$$

where  $U \in \mathbb{R}^{m \times r}$  has as its  $i$ th row the vector  $u_i^\top$ , and similarly for  $V$ . The key aspect of this model is that  $U$  and  $V$  do not depend on  $k$ ; that is, the latent features are the same across matrices. This assumption relates the various matrices  $M^k$  to each other and allows for the potential of using other slices for learning. In particular, in addition to  $X^1$ , data from  $X^2, \dots, X^n$  can be used to improve the rate of recovery on  $M^1$ . Note that the possibility for some elements of  $S^k$  to be equal to zero means that different interactions do not need to involve the *same* latent factors, but rather all interactions draw from a small, shared pool of latent factors. As before, the generality of this model is determined by the dimension of the latent space,  $r$ , a quantity that we will formalize as the *slice rank* of the tensor  $M$ .

**Definition 1.** The slice rank of a tensor  $M \in \mathbb{R}^{m \times m \times n}$ , denoted as  $\text{Slice}(M)$ , is the maximum of the ranks of its mode 1 and mode 2 unfoldings; that is,  $\text{Slice}(M) = \max(\text{rank}(M_{(1)}), \text{rank}(M_{(2)}))$ .

In the definition above,  $\text{rank}(M_{(1)})$  is the number of latent customer features and  $\text{rank}(M_{(2)})$  is the number of latent product features. Because these numbers may not be equal (corresponding to having a different number of customer and product features), the slice rank is defined as the maximum of the two quantities.

In what follows, we will seek to recover tensors of low *slice rank* from their noisy observations—a formalization of the estimation problem that motivates this paper. Before doing so, we find it worthwhile to note that there is a long history of using tensors as a meaningful data structure to capture multivariate data, and this literature has yielded other notions of

tensor rank. Specifically, the two most common of these notions are the *canonical (or CP) rank* and the *Tucker rank*, which have both been used in applications as diverse as psychometrics (Carroll and Chang 1970), chemical analysis (Henrion 1994), facial recognition (Vasilescu and Terzopoulos 2002), chatroom analysis (Acar et al. 2005), and web search (Sun et al. 2005). As it turns out, the requirement of a low slice rank is *weaker* than requiring either low CP rank or low Tucker rank. Specifically, denote by  $CP(r)$  the set of tensors  $M \in \mathbb{R}^{m \times m \times n}$  that have a CP rank at most  $r$ . Similarly, denote by  $\text{Tucker}(r, r, l)$  the set of tensors  $M \in \mathbb{R}^{m \times m \times n}$  with a Tucker rank at most  $(r, r, l)$  (see Section B of the online appendix for a formal definition of CP rank and Tucker rank). We then have the following.

**Proposition 1.** *The set of tensors  $M \in \mathbb{R}^{m \times m \times n}$  with a slice rank at most  $r$ ,  $\text{Slice}(r)$ , contains  $CP(r)$  and  $\text{Tucker}(r, r, l)$ :*

$$\mathcal{CP}(r) \subseteq \text{Slice}(r), \quad \text{and} \quad \text{Tucker}(r, r, l) \subseteq \text{Slice}(r).$$

The proof of this result can be found in Section B of the online appendix. In summary, this result establishes that our requirement of “simplicity” subsumes important existing notions of simplicity for tensors. More importantly, as is evident from our presentation thus far, the notion of low slice rank has an elegant interpretation in our context as seeking out latent representations for customers and products along with functional forms that relate these latent representations to the likelihood of a specific interaction.

### 2.3. Our Problem

Our problem is to recover the tensor  $M$  from a noisy observation of its entries. In particular, our observation consists of the data tensor  $X = M + \epsilon$ , where the elements of the “noise” tensor  $\epsilon$  are independent with mean zero. We emphasize that, so far, we have not restricted the elements of  $\epsilon$  to be Gaussian or even identically distributed, as is typically done.<sup>2</sup>

Our goal is to construct an estimator  $\hat{M}(X)$  to minimize some loss function with respect to  $M$ . To reduce notation, we will use  $\hat{M}(X)$  and  $\hat{M}$  interchangeably. The usual loss function that we will take here is the mean-squared error (MSE):

$$\text{MSE}(\hat{M}) = \mathbb{E} \left[ \sum_{k=1}^n \frac{1}{nm^2} \|\hat{M}^k - M^k\|_F^2 \right].$$

We will refer to the problem of constructing an estimator to minimize MSE as *tensor recovery*. In addition, as we have noted, in many cases we might only be interested in recovering a single slice of the tensor (having observed all of  $X$ ). For example, even with data from many types of customer–product interactions, we may be solely interested in predicting purchase probabilities. In these settings, the MSE is not an ideal

loss function, as it measures average recovery error over all slices. Therefore, in addition to MSE, we will also consider the *slice mean-squared error* (SMSE):

$$\text{SMSE}(\hat{M}) = \max_{1 \leq k \leq n} \mathbb{E} \left[ \frac{1}{m^2} \|\hat{M}^k - M^k\|_F^2 \right],$$

and refer to the problem of constructing an estimator to minimize SMSE as *slice recovery*. SMSE is a more robust loss function, as it measures the maximum recovery error over all slices, and so a guarantee on the SMSE applies to every single slice. Also, because SMSE is always greater than or equal to MSE, any upper bound for SMSE will apply to MSE—and therefore the tensor recovery problem.

### 3. A Lower Bound and Incumbent Approaches

In the previous section, we recalled a minimax lower bound on the error rate one may hope to achieve under any algorithm for the problem of matrix recovery. Indeed, this bound motivated our problem of tensor recovery, where, loosely stated, we hope to use data from multiple types of customer–product interactions to improve our estimate of the likelihood of a given type of interaction. Our goal in this section will be to understand the extent to which data on additional types of interactions can help with the recovery of a specific type of interaction. Thus motivated, we next establish a lower bound on the error rate that one may hope to achieve under *any* algorithm for our problem.

**Proposition 2.** *For any estimator  $\hat{M}$ , we can construct a tensor  $M \in \text{Slice}(r)$  with entries in  $[-1, 1]$  and a random noise  $\epsilon$  with independent, zero mean entries, such that  $X = M + \epsilon$  has entries in  $[-1, 1]$  almost surely, and*

$$\text{SMSE}(\hat{M}) \geq \text{MSE}(\hat{M}) \geq C \left( \frac{r^2}{m^2} + \frac{r}{nm} \right),$$

where  $C$  is a universal constant.

The proof of Proposition 2 is presented in Section C of the online appendix, and it relies on a carefully constructed ensemble of problem instances. The proposition lets us draw several key conclusions with regard to the power of side information.

1. *The special case of matrix recovery:* In the matrix recovery setting—that is,  $n = 1$  so that  $M$  is a matrix of rank  $r$ —we recover a minimax lower bound of  $\text{MSE}(\hat{M}) = \Omega(r/m)$ . This bound is well known (e.g., Candès and Plan 2011), and a number of matrix recovery algorithms achieve this bound. Consequently, the naïve approach of using an optimal matrix recovery algorithm *separately on each slice* achieves  $\text{MSE} = O(r/m)$  and  $\text{SMSE} = O(r/m)$ . Beating such an approach is precisely the motivation for our work here.

2. *The potential benefit of side information:* The impact of side information is made precise by the dependence, on  $n$ , of our lower bound on achievable error rate. The minimax bound above suggests that additional side information might permit up to a *linear* reduction in recovery error up to a certain point, beyond which additional side information cannot help. Specifically, the lower bound is dominated by an error term that scales like  $r/nm$  for  $n$  sufficiently small, and it can be no smaller than  $r^2/m^2$  irrespective of the amount of side information  $n$ . *How does side information alleviate the sparsity problem?* Recall that from our motivating example in the introduction, the minimax optimal recovery rate of  $O(r/m)$  achievable with data of a single type ( $n = 1$ ) implies that we need to observe on the order of  $r$  purchases per user, which is a problem because users make on average less than a single purchase a month. In the setting above, we could get by with on the order of just  $r/n$  purchases per user as long as we recover a total of  $r^2$  purchases across all users and products. Put another way, a retailer that can easily collect data on  $10^2$  types of interactions can hope to estimate a much richer model with  $r = 10^2$ .

In summary, the minimax lower bound established above raises the specter of dramatically increasing our ability to cope with sparse data, provided we have access to sufficient side information. Do existing algorithms come close to achieving this lower bound?

#### 3.1. Incumbent Approaches

The dominant approach to tensor recovery relies on convex optimization. To motivate this approach by example, consider the recovery of a low Tucker rank tensor from noisy observation  $X$ . We can formulate this problem as the following (hard) optimization problem:

$$\begin{aligned} \min_Y \quad & \|Y - X\|_F^2 \\ \text{s.t.} \quad & \text{rank}(Y_{(i)}) \leq r_i, \quad i = 1, 2, 3. \end{aligned} \tag{1}$$

That is, we must choose a tensor  $Y$  that most closely “matches” the observed data  $X$  from the set of tensors with Tucker rank bounded by  $(r_1, r_2, r_3)$ . For certain noise models such as, for example, independent and identically distributed (i.i.d.) Gaussian noise, the solution to (1) would correspond to the maximum likelihood estimator.

Because this is a difficult problem (as a result of the rank constraints), the standard approach taken is to come up with a convex surrogate for the tensor rank. So, for example, one such variant is

$$\min_Y \|X - Y\|_F^2 + \sum_{i=1}^3 \lambda_i \|Y_{(i)}\|_{\ast}, \tag{2}$$

where the ranks are replaced by the nuclear norms of the appropriate unfoldings, and in addition, the rank constraints have also been dualized (in keeping with how such relaxations are presented in the literature); the



weights  $\lambda_i$  are chosen by the user and intuitively should encode prior knowledge of rank. This convex algorithm has been studied extensively (Gandy et al. 2011, Tomioka et al. 2011, Liu et al. 2013, Signoretto et al. 2014). In the case where  $n = 1$ , the mode 3 unfolding would naturally be removed from the objective, in which case the algorithm is precisely equivalent to the de facto convex formulation for matrix recovery, rendering the formulation a natural one. Tomioka et al. (2011) show that if the noise tensor has i.i.d. Gaussian entries with standard deviation  $\sigma$ , the estimator above achieves

$$\text{MSE}(\hat{M}) = O(\sigma^2(r/m + r/n)) \quad (3)$$

if the Tucker rank of  $M$  is  $(r, r, r)$ . Recall that any tensor with a Tucker rank  $(r, r, r)$  has a slice rank at most  $r$ , so this result applies to a subset of tensors with a slice rank  $r$ . There is good reason to believe that this guarantee is tight. For example, Tomioka et al. (2011) also show a very similar guarantee in a related setting where random linear combinations of the entries of  $M$  are observed, and Mu et al. (2013) show that, in fact, that guarantee is tight. Furthermore, in our experiments later on, we will empirically observe that the recovery rate scales as (3).

Taken together, this is disappointing—it shows that the convex approach above does not improve on recovering individual slices of the matrix via an optimal matrix completion algorithm and leaves a wide chasm between the minimax lower bound of Proposition 2 and the error rate achieved. Put simply, this approach *does not solve the data sparsity problem*. The issue of why existing convex approaches do not achieve optimality is a very interesting question. The main reason to hope that convex approaches might work in the tensor setting is because of their success in matrix recovery where the nuclear norm is easily shown as the *tightest* convex relaxation of the rank function. The major challenge in going from matrices to tensors is that there is not an obvious “optimal” convexification of (1) above. Formulation (2) is one possible convexification, but it is provably *not* a tight convexification. Other convex problems have been suggested—for instance, in Romera-Paredes and Pontil (2013), Mu et al. (2013), and Tomioka and Suzuki (2013)—but none of these proposals improves on the recovery guarantee above.

Outside of convex formulations, Suzuki (2015) recently proposed a Bayesian estimator that matches the lower bound in Proposition 2 for i.i.d. Gaussian noise and tensors with low CP rank (as is evident from the definition of CP rank, this is significantly more restrictive than slice rank). Unfortunately, this procedure relies on a Monte Carlo approach in high dimension, and its computational efficiency is unknown (i.e., we no longer have the computational efficiency guarantees that come with the convex approach).

In summary, we may conclude that employing existing tensor recovery machinery for the problem at hand does not yield a solution to the data sparsity problem. In fact, using the de facto convex approaches for the problem cannot be expected to yield any improvement over the approach of individually recovering each slice of data. Now, in contrast to taking the convex relaxation approach above, our algorithm works by directly (and efficiently) constructing a *feasible* solution to the original problem (1). The constructed solution is not necessarily an optimal solution to (1), but our core theoretical result shows that it is minimax optimal for noisy recovery.

#### 4. An Algorithm for Slice Recovery

We will now present our algorithm for slice recovery and tensor recovery, which we will see is of a fundamentally different nature than the convex approaches, as just described. Recall from the setup that  $X$  is our observed data tensor and  $M$  is the ground truth tensor we are trying to recover. We assume  $M$  has slice rank  $r$ , which implies the existence of a decomposition wherein every slice  $M^k$  can be represented as  $M^k = US^kV^T$ , where  $U, V$  are  $m \times r$  matrices encoding latent customer and product features, and each  $S^k$  is an  $r \times r$  matrix that captures the specific bilinear form for each slice  $M^k$ .

Up to this point, it has been convenient to think of each column of  $U$  and  $V$  as encoding a specific, possibly interpretable feature (e.g., customer demographics, product specifications), but in fact, our algorithm is best understood when  $U$  and  $V$  are viewed as latent feature spaces. In particular,  $U$  and  $V$  each encode a linear subspace of  $\mathbb{R}^m$ , the subspaces spanned by their respective columns. Note that because our model places no restriction on the bilinear interaction terms  $S^k$ , the terms  $U$  and  $V$  in the decomposition  $M^k = US^kV^T$  are only unique up to the subspaces they span; that is, we could replace  $U$  and  $V$  with any set of features that span the same feature space. For this reason, we will refer to  $U$  and  $V$  as features and feature spaces interchangeably, and for mathematical convenience, we will assume without loss of generality that the columns of  $U$  and  $V$  are orthonormal.

The algorithm proceeds in two stages, both of which we motivate from first principles.

**Stage 1: Learning Subspaces.** In the first stage, we use data from every slice to estimate the latent feature spaces  $U$  and  $V$ . Let us first focus on the procedure for learning  $U$ . Our first step in this regard is to construct the mode 1 unfolding  $X_{(1)}$ , which, recall, is the  $m \times nm$  matrix whose columns are the columns of  $X^1, \dots, X^n$  (the order of the columns does not matter). We then compute  $\hat{U}$ , our estimate of  $U$ , as the first  $r$  left-singular vectors of  $X_{(1)}$ . More precisely, assuming that  $X_{(1)}$  admits the singular value decomposition  $X_{(1)} = U_1 \Sigma_1 V_1^T$ , we set  $\hat{U}$  to be the

columns of  $U_1$  corresponding to the  $r$  largest singular values. We denote this entire procedure with the shorthand

$$\hat{U} = \text{svds}(X_{(1)}, r). \tag{4}$$

Under very mild assumptions, it will turn out that  $\hat{U}$  is a “good” estimate of  $U$ . To see this, we can view  $X_{(1)}$  as a noisy observation of  $M_{(1)}$ , which is a wide matrix with a row for each customer and a column for each product–interaction type pair. The key point is that the columns of  $M_{(1)}$  span the feature space  $U$ , so if we were allowed to observe  $M_{(1)}$ , we could easily find  $U$  as the space spanned by the columns of  $M_{(1)}$  (i.e., its “column space”). Instead, we observe  $X_{(1)}$ , which, because of the added noise, does not have column space  $U$ ; in fact, its columns are likely to span all of  $\mathbb{R}^m$ . Still, under mild assumptions on the nature of the noise, we can expect that the columns of  $X_{(1)}$  lie approximately in  $U$  and therefore estimate  $\hat{U}$  as the orthonormal matrix that minimizes the function  $f$  defined as

$$f(U_1) = \min_{R_1 \in \mathbb{R}^{r \times nm}} \|U_1 R_1 - X_{(1)}\|_F.$$

For any  $r$ -dimensional subspace  $U_1$ ,  $f(U_1)$  measures how closely  $X_{(1)}$  can be approximated by a matrix with column space  $U_1$ , and its minimizer  $\hat{U}$  is precisely (4). This stage is exactly where we take advantage of having multiple slices of data, as the accuracy of  $\hat{U}$  as an estimate of  $U$  is better the more slices we have, or the wider  $M_{(1)}$  is. We will quantify this exactly in the next section.

To estimate  $V$ , we apply a similar procedure using the mode 2 unfolding  $X_{(2)}$ , which, recall, is the  $m \times nm$  matrix whose columns are the transposed rows of  $X^1, \dots, X^n$ . Just as in the previous discussion,  $X_{(2)}$  is a noisy observation of  $M_{(2)}$ , whose column space is  $V$ . It follows that, under some assumptions on the noise, the columns of  $X_{(2)}$  lie approximately in  $V$ , and a natural estimate of  $V$  is

$$\hat{V} = \text{svds}(X_{(2)}, r).$$

**Stage 2: Projection.** The second stage works on each slice separately. Having estimated  $U$  and  $V$  as  $\hat{U}$  and  $\hat{V}$ , respectively, it remains to estimate the bilinear terms  $S^k$  for each slice  $M^k$ . We do this by solving the following optimization problem for each slice:

$$\hat{S}^k = \arg \min_S \|\hat{U} S \hat{V}^\top - X^k\|_F. \tag{5}$$

To motivate this, suppose that instead of  $\hat{U}$  and  $\hat{V}$ , we had access to  $U$  and  $V$  exactly. Then knowing that  $M^k$  takes the form  $US^kV^\top$  for some  $S^k$ , our estimate of  $M^k$  should take this same form. Therefore, we would use the closest approximation to  $X^k$  of this form, essentially

“projecting”  $X^k$  onto the feature spaces  $U$  and  $V$ . Because we have only  $\hat{U}$  and  $\hat{V}$ , we use the closest approximation of the form  $\hat{U} S^k \hat{V}^\top$  instead, as in (5).

The above is a least-squares problem and, as such, admits a closed-form solution: assuming that  $\hat{U}$  and  $\hat{V}$  are orthonormal, the solution of (5) is  $\hat{S}^k = \hat{U}^\top X^k \hat{V}$ , and so our final estimate of  $M^k$  can be written as

$$\hat{M}^k = \hat{U} \hat{U}^\top X^k \hat{V} \hat{V}^\top.$$

One nice interpretation here is that  $\hat{U} \hat{U}^\top$  and  $\hat{V} \hat{V}^\top$ , which we will denote by  $P_{\hat{U}}$  and  $P_{\hat{V}}$ , respectively, are projection operators: left-multiplying a matrix by  $P_{\hat{U}}$  replaces each of its columns with the orthogonal projection onto the space  $\hat{U}$ , and similarly, right-multiplying by  $P_{\hat{V}}$  replaces rows by the orthogonal projection onto  $\hat{V}$ . The resulting matrix lies in the feature spaces  $\hat{U}$  and  $\hat{V}$ .

The entire slice learning algorithm is outlined in Algorithm 1. The algorithm and our forthcoming analysis are easily extended to using different values of  $r$  when forming  $\hat{U}$  and  $\hat{V}$ . This may be advantageous when the latent customer and product feature spaces of  $M$  are of significantly different dimensions, but we use a single value  $r$  here for ease of notation.

**Algorithm 1** (Slice Learning)

Input:  $X, r$

1.  $\hat{U} \leftarrow \text{svds}(X_{(1)}, r)$
2.  $\hat{V} \leftarrow \text{svds}(X_{(2)}, r)$
3.  $\hat{M}^k \leftarrow P_{\hat{U}} X^k P_{\hat{V}}, \quad k = 1, \dots, n$

Output:  $\hat{M}^k, k = 1, \dots, n$

**4.1. Practical Considerations**

We conclude this section with a discussion of practical implementation and computation issues.

**4.1.1. Knowing  $r$ .** The slice learning algorithm above takes the slice rank  $r$  as input, but in some settings we may not know  $r$  in advance. In particular, we do not know the ranks of the two unfoldings  $M_{(1)}$  and  $M_{(2)}$  in advance. This is a common challenge in low-rank matrix recovery, although in that setting, the problem has proven to be relatively benign. Specifically, there exists a wide array of rank estimation methods that we may borrow from, including cross validation (Wold 1978, Owen and Perry 2009), visual inspection of plotted singular values (Cattell 1966), and Bayesian methods (Hoff 2007).

Perhaps the simplest approach when  $r$  is not known is a “universal” thresholding scheme where we preserve only singular vectors corresponding to singular values above a certain easily precalculated threshold. This has been shown to work in the matrix recovery setting (Chatterjee 2014, Gavish and Donoho 2014), and we anticipate that such a scheme will work just as

effectively here. In particular,  $X_{(1)}$  is the sum of the signal  $M_{(1)}$  and the noise  $\epsilon_{(1)}$ , so if the singular values of  $\epsilon_{(1)}$  are all significantly lower than the  $r$  nonzero singular values of  $M_{(1)}$ , then by Weyl’s inequality,  $X_{(1)}$  will have  $r$  singular values that are much larger than the rest.

To make this more precise, consider the following argument, which also serves as an introduction to the type of arguments used in the next section. Suppose that the terms of  $\epsilon_{(1)}$  are i.i.d. with unit variance and bounded fourth moment; we will assume a more general noise model later. To fix the signal-to-noise ratio, assume that the terms of  $M_{(1)}$  are of constant order, so  $\|M_{(1)}\|_F^2 = \Theta(m^2n)$ . First,  $M_{(1)}$  has rank  $r$  (for some unknown  $r$ ), so if the nonzero singular values of  $M_{(1)}$  are all within a constant order of each other, then these singular values will be of  $\Theta(\|M_{(1)}\|_F/\sqrt{r})$ , or  $\Theta(m\sqrt{n/r})$ . At the same time, by the Bai-Yin law (Yin et al. 1988), the largest singular value of  $\epsilon_{(1)}$  does not scale greater than  $O(\sqrt{mn})$ . Therefore, a clear separation forms asymptotically between the nonzero singular values of  $M_{(1)}$  and the singular values of  $\epsilon_{(1)}$ . By Weyl’s inequality,  $X_{(1)}$  will have  $r$  singular values of  $\Theta(m\sqrt{n/r})$ , and its remaining singular values are of  $O(\sqrt{mn})$ . It follows that choosing a threshold in this gap and retaining the singular values greater than this threshold will closely approximate using the correct value of  $r$ . This result can be formalized as a theorem, but we do not do so here.

**4.1.2. Computation with Large Tensors.** An important consideration is the scale at which the slice learning algorithm needs to operate—the nature of the applications necessitates that predictions be made rapidly, as the data are constantly changing and up-to-date output is required. Therefore, computational efficiency is of paramount importance in practice.

The only computationally intensive step in the slice learning algorithm is in Stage 1, which requires the computation of two partial SVDs. Again, the largest retailers have upwards of  $10^8$  unique customers and products ( $m \sim 10^8$ ), and they can easily collect data on hundreds of interactions ( $n \sim 100$ ), at the very least. For dense matrices, this massive scale renders generic SVD algorithms intractable. Fortunately, although the ambient dimensions of the input data are large, the data are itself typically quite sparse. Data sparsity has so far been treated as a disadvantage, because it limits the complexity of the models we can learn, but it is a key advantage from a computational standpoint. There exist mature linear algebraic algorithms that compute the top singular vectors while exploiting data sparsity—for example, using Lanczos iterations as in the PROPACK algorithm (Larsen 1998). Specifically, these algorithms rely on a power method that repeatedly applies a

matrix-vector multiplication subroutine. Because the matrix in question (the unfolding of the sparse tensor) is sparse, this operation can be implemented with running time linear in the number of nonzero elements of the matrix. This is already drastically less computation than convex approaches to this task, which by and large are solved with iterative algorithms (Gandy et al. 2011) that require performing dense SVDs multiple times.

**4.1.3. Customer/Product Side Information.** Our generative model includes a set of latent customer and product features  $U$  and  $V$ , and Stage 1 of our algorithm essentially works by estimating these features jointly across all slices of data. We will see that the performance of our algorithm boils down to how well  $U$  and  $V$  can be estimated, and that going from a single slice of data to many slices interpolates between standard matrix recovery (with no prior knowledge of  $U$  or  $V$ ) and recovery given  $U$  and  $V$  exactly.

In many cases, we may have side information in the form of explicit features about customers or products that are believed to be relevant (e.g., customer demographics, product specifications). This is the subject of a line of work in matrix recovery (Xu et al. 2013, Soni et al. 2016) that seeks to recover a single slice, assuming that the feature spaces  $U$  and  $V$  are known exactly.

The approach we have laid out is flexible enough to incorporate the same kind of customer and product features without requiring the assumption of knowing  $U$  and  $V$  exactly. Suppose that, in addition to the data tensor  $X$ , we have  $\ell$  customer features that we encode in the  $m \times \ell$  matrix  $A$ . We can include this information in our algorithm by expanding the estimated feature space  $\hat{U}$  with these extra features; that is, after producing the estimate  $\hat{U}$  in Stage 1 as usual, we can replace  $\hat{U}$  with the subspace spanned by the columns of  $\hat{U}$  and  $A$  and then execute Stage 2 normally. The equivalent procedure can be done with  $\hat{V}$  if we have product features.

## 5. Recovery Guarantees for Slice Learning

The goal of this section is to provide a statistical recovery guarantee for our slice learning algorithm. In particular, we are looking for a guarantee that improves on the naïve approach of ignoring side information and recovering each slice separately. As discussed in Section 3, that naïve approach, which is effectively a matrix recovery algorithm, can achieve  $\text{SMSE} = O(r/m)$ .

So far, we have made the assumption that the ground truth tensor has low slice rank, which as we have seen, restricts the complexity of the underlying generative model and offers the possibility of improving on the matrix recovery rate. However, beating this rate is impossible without making further assumptions. Specifically, consider the case where  $M$  has  $n - 1$  slices,

all of whose entries are identically 0, and a single nontrivial, low-rank slice, which we take without loss of generality as the first slice. Furthermore, assume that the noise tensor has independent unit variance Gaussian entries on the first slice and is identically zero on the remaining  $n - 1$  slices. Although the ground truth tensor has low slice rank, the problem of recovering the first slice is now literally no different from the problem of matrix recovery on the first slice, because the remaining slices are superfluous.

To this end, we define a structural parameter for  $M$  that we will eventually see controls the rate at which learning across slices is possible. Letting  $\sigma_r^2(M_{(1)})$  and  $\sigma_r^2(M_{(2)})$  denote the  $r$ th-largest singular values of  $M_{(1)}$  and  $M_{(2)}$ , respectively, we define the *learning rate* as follows.

**Definition 2.** The learning rate of a tensor  $M \in \mathbb{R}^{m \times m \times n}$  with slice rank  $r$ , denoted by  $\gamma_M$ , is defined as

$$\gamma_M = \frac{r}{m^2 n} \min(\sigma_r^2(M_{(1)}), \sigma_r^2(M_{(2)})).$$

We will shortly see how  $\gamma_M$  plays the role of a rate of learning. For now, we merely comment on the range of values one might reasonably expect this quantity to take. On the low end, observe that, by Marčenko and Pastur (1967), if the noise tensor were i.i.d., the (squared) singular values of the noise tensor unfoldings  $\epsilon_{(1)}$  and  $\epsilon_{(2)}$  are  $O(mn)$ . A minimal requirement is that the singular values of  $M_{(1)}$  and  $M_{(2)}$  dominate those of the noise unfoldings, which would, in turn, imply that  $\gamma_M = \Omega(r/m)$ , so the loosest requirement we can reasonably place on  $\gamma_M$  is to require  $\gamma_M = \Omega(r/m)$ .

At the other end of the spectrum, given that the entries of  $M$  are required to be bounded, we must have that  $\|M_{(1)}\|_F^2 = O(m^2 n)$ , so that  $\sigma_r^2(M_{(1)}) = O(m^2 n/r)$ . This, in turn, implies that  $\gamma_M = O(1)$ , and the strongest requirement we can place is to require  $\gamma_M$  be a constant. In fact, in terms of scaling, it is reasonable to treat  $\gamma_M$  as constant: the condition  $\|M_{(1)}\|_F^2 = \Theta(m^2 n)$  is satisfied for all but trivial examples such as the one just described, and then  $\sigma_r^2(M_{(1)}) = \Theta(m^2 n/r)$  follows as long as the largest and smallest nonzero singular values of  $M_{(1)}$  are not significantly different; that is, a parameter akin to the condition number is bounded.<sup>3</sup>

Before proceeding with a statement of our main result, we will place an assumption on the noise (our only nontrivial assumption thus far). Specifically, we will require the noise to be “balanced” in a certain sense.

**Assumption 1 (Balanced Noise).** Let  $v$  be the tensor whose entries are the variances of the corresponding entries of  $\epsilon$ . Specifically,  $v_{i,j}^k = E[(\epsilon_{i,j}^k)^2]$ . The noise  $\epsilon$  is said to be balanced if the row-sums of  $v_{(1)}$  are all equal and the row-sums of  $v_{(2)}$  are all equal.

The assumption above is trivially satisfied in the case of i.i.d. noise, and, in particular, the case of i.i.d. Gaussian noise that is often studied in the matrix and tensor recovery setting. Refinements of our result allow for weaker versions of this assumption; we will see in the next subsection that we can permit a certain amount of discrepancy in the row sums of  $v_{(1)}$  and  $v_{(2)}$ , and we allow this to grow with  $m$  and  $n$ . We are now ready to state our main result.

**Theorem 1 (Balanced Noise).** Assume that the entries of  $M$  lie in  $[-1, 1]$ . If the entries of  $\epsilon$  are independent, mean-zero, and  $E[(\epsilon_{ij}^k)^6] \leq K^6$ , and if, furthermore,  $\epsilon$  is balanced, then there exists a constant  $c(K)$  that depends only on  $K$  such that for the slice learning algorithm,

$$\text{MSE}(\hat{M}) \leq \text{SMSE}(\hat{M}) \leq c(K) \left[ \frac{r^2}{m^2} + \frac{r^2}{\gamma_M^2 mn} \right].$$

The proof of Theorem 1 can be found in Section D in the online appendix. We next evaluate this result in light of the minimax guarantee established in Proposition 2 and, more generally, our broader goal of slice recovery:

1. *Learning from slices:* As discussed earlier, at the very least we expect  $\gamma_M = \Omega(r/m)$ . Even in this setting, we see that provided that  $n$  is sufficiently large, Theorem 1 guarantees an SMSE (and, consequently, an MSE as well) that is  $O(r^2/m^2)$ . In particular, for  $n$  sufficiently large, we obtain a recovery rate that meets the leading term of the minimax guarantee in Proposition 2. Of course, this is *substantially* better than the available guarantee for the naïve approach, which was  $O(r/m)$ .

2. *High learning rate:* As  $\gamma_M$  grows, so too does our ability to learn across slices. Specifically, for  $\gamma_M = \Theta(1)$ , Theorem 1 guarantees  $\text{MSE} \leq \text{SMSE} = O(r^2/m^2 + r^2/mn)$ . This is within a factor of  $r$  off from the lower bound of  $\text{MSE} = \Omega(r^2/m^2 + r/mn)$  in Proposition 2. Put a different way, we achieve  $\text{SMSE} = O(r^2/m^2)$  with only  $n = \Omega(m)$  slices of side information.

3. *Recovery rate improvement:* The recovery rate in Theorem 1 is for low slice rank tensors, which includes tensors with low Tucker rank and CP rank. We emphasize that the rate  $\text{MSE} = O(r^2/m^2 + r^2/mn)$  greatly improves on the best-known theoretical guarantees for noisy recovery of low Tucker rank tensors and for convex optimization approaches to recovering low CP rank tensors.

### 5.1. Relaxing the Balanced Noise Assumption

Although the balanced noise assumption is already a generalization of many frequently studied noise models, it is worth considering how our algorithm performs when this assumption does not hold. To do so, we will present a more general version of Theorem 1 that relaxes the balanced noise assumption and reflects the recovery error caused by “unbalanced” noise.

To proceed, we need to precisely quantify the concept of *unbalanced* noise. Recall that if  $v$  is the tensor whose entries are the variances of the corresponding entries of  $\epsilon$  (i.e.,  $v_{ij}^k = E[(\epsilon_{ij}^k)^2]$ ), then  $\epsilon$  is balanced if the row-sums of  $v_{(1)}$  are equal and the row-sums of  $v_{(2)}$  are equal. An equivalent way to state this assumption is  $E[\epsilon_{(1)}\epsilon_{(1)}^\top] = \rho_1 I_m$  and  $E[\epsilon_{(2)}\epsilon_{(2)}^\top] = \rho_2 I_m$  for some constants  $\rho_1$  and  $\rho_2$ , where  $I_m$  is the  $m \times m$  identity matrix. To see this, note that the off-diagonal elements of  $E[\epsilon_{(1)}\epsilon_{(1)}^\top]$  and  $E[\epsilon_{(2)}\epsilon_{(2)}^\top]$  are always equal to zero when the noise terms are independent, and the diagonal elements are exactly the row-sums of  $v_{(1)}$  and  $v_{(2)}$ , so the balanced noise assumption states that  $E[\epsilon_{(1)}\epsilon_{(1)}^\top]$  and  $E[\epsilon_{(2)}\epsilon_{(2)}^\top]$  are multiples of the identity matrix.

For general, possibly unbalanced noise, it turns out that the appropriate quantities to measure the level of “unbalance” in the noise are

$$\min_{\rho} \frac{1}{m} \left\| E[\epsilon_{(1)}\epsilon_{(1)}^\top] - \rho I_m \right\|_F^2 \quad \text{and} \\ \min_{\rho} \frac{1}{m} \left\| E[\epsilon_{(2)}\epsilon_{(2)}^\top] - \rho I_m \right\|_F^2.$$

These quantities measure how far  $E[\epsilon_{(1)}\epsilon_{(1)}^\top]$  and  $E[\epsilon_{(2)}\epsilon_{(2)}^\top]$  are from a multiple of the identity matrix. One nice interpretation is that the quantities correspond to population variances, one each for the row-sums of  $v_{(1)}$  and the row-sums of  $v_{(2)}$ . We denote the maximum of these two quantities as  $\delta^2$ , and can now state our more general result.

**Theorem 2.** *Assume the entries of  $M$  lie in  $[-1, 1]$ . If the entries of  $\epsilon$  are independent, mean-zero, and  $E[(\epsilon_{ij}^k)^6] \leq K^6$ , then there exists a constant  $c(K)$  that depends only on  $K$  such that for the slice learning algorithm,*

$$\text{SMSE}(\hat{M}) \leq c(K) \left[ \frac{r^2}{m^2} + \frac{r^2}{\gamma_M^2 mn} + \frac{r^2 \delta^2}{\gamma_M^2 m^3 n^2} \right].$$

To interpret the guarantee in Theorem 2, consider the range of values that  $\delta^2$  might take. The lowest possible value is  $\delta^2 = 0$ , which corresponds to the balanced noise setting, and in this case, we recover Theorem 1 exactly. On the other hand, each row of  $v_{(1)}$  and  $v_{(2)}$  contains  $mn$  elements, so their row-sums may scale as  $O(mn)$ , and thus in the worst case,  $\delta^2 = O(m^2 n^2)$ . In this worst-case setting, Theorem 2 does not improve on the guarantee for the naïve approach. However, the recovery rate in Theorem 2 matches that in Theorem 1 as long as  $\delta^2 = O(m^2 n)$ . Also, because  $\delta^2$  measures the population variances of the two sets of row-sums, it is highly robust to settings where only a few row-sums are significantly unbalanced. All of this suggests that Theorem 1 holds even if the balanced noise assumption is significantly relaxed.

To this point, our work has applied to the problem of noisy tensor recovery, a framework that addresses

settings such as the retail example and our experiment with music streaming data. As discussed in Section 1, there are settings and applications where the existing data instead can be represented as a partially observed tensor (i.e., the tensor completion problem). The challenge here is to design algorithms that are optimal in terms of the *number of observed entries* required for exact recovery. Now, under the assumption that entries are observed uniformly at random, it is possible to map completion problems to noisy recovery problems using a technique developed for matrix completion (Achlioptas and McSherry 2007, Keshavan et al. 2010, Chatterjee 2014). We discuss this topic in Section G of the online appendix, where we (a) use this same device to adapt the slice learning algorithm to tensor completion, (b) state a corollary to Theorem 2 that characterizes the requisite number of observed entries, (c) compare this result to existing tensor completion algorithms, and (d) show results of synthetic tensor completion experiments.

## 6. Experiments

We performed two sets of experiments to evaluate the slice learning algorithm, the first using randomly generated tensors and the second using a real-world dataset. This section describes these experiments and their results in detail, wherein the following points emerge:

1. The slice learning algorithm drastically outperforms convex algorithms by an order of magnitude, even though convex algorithms require drastically greater computation. On tensors with low slice rank, the slice learning algorithm outperforms a convex benchmark we propose that is in the spirit of incumbent approaches; this is a “home court” setting, as the recovery guarantees from the previous section show that our algorithm is expected to recover these tensors. More surprisingly, on tensors with low Tucker rank, the slice learning algorithm also outperforms an existing convex benchmark designed specifically for low Tucker rank tensors.

2. On real-world, sparse data, the slice learning algorithm outperforms two common benchmark approaches. This suggests that real data, when represented as a tensor, indeed exhibit the type of structure that the slice learning algorithm is able to exploit.

3. By leveraging sparsity, the slice learning algorithm can be used on large data sets, in regimes where common operations such as taking the SVD of a dense matrix are intractable.

### 6.1. Synthetic Recovery Experiments

We conducted a number of experiments on randomly generated tensors to test the performance of the slice learning algorithm, for learning individual slices and the tensor as a whole. In all of our experiments, we compare against a benchmark convex optimization approach.

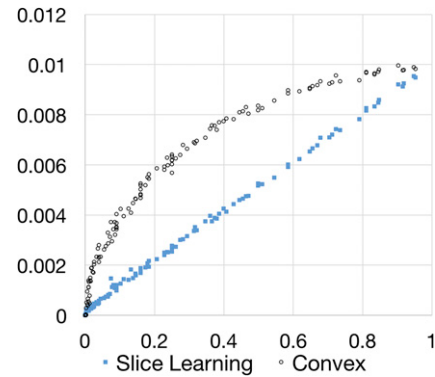
**6.1.1. Noisy Recovery and Tucker Rank.** For our first experiment, we randomly generated  $m \times m \times m$  tensors with Tucker rank  $(r, r, r)$ , where in each replication,  $m$  was drawn uniformly from the integers between 10 and 50, and  $r$  was then drawn uniformly from the integers between 1 and  $m$ . In each replication, we randomly drew orthonormal  $m \times r$  matrices  $U, V$ , and  $W$ , along with an  $r \times r \times r$  tensor  $S$  with entries drawn from the standard normal distribution. The ground truth tensor was then constructed in the canonical way from  $U, V, W$ , and  $S$  (see Section B in the online appendix), and the observed data tensor  $X$  was constructed by adding independent mean-zero Gaussian noise of standard deviation 0.1 to each entry. This is nearly identical to the experimental setup in Tomioka et al. (2011).

We compared the slice learning algorithm against the convex algorithm that minimizes (2) in Section 3.1. Recall that this is one of the few well-studied algorithms with a theoretical recovery guarantee (Tomioka et al. 2011). We solved this using the Douglas–Rachford splitting method, as described by Gandy et al. (2011). Note that the slice learning algorithm requires an estimated rank as input, and in this experiment, the algorithm was given the true rank  $r$  in each replication. On the other hand, the convex objective (2) has a parameter  $\lambda$  that encodes knowledge of the rank; to level the playing field, in each iteration we solved (2) for values of  $\lambda$  ranging from  $2^{-2}$  to  $2^5$  and reported the best performance among all of these.

We performed 100 replications. The results are depicted in Figure 4, where each replication is represented by two points, one for each algorithm. Each point represents the MSE versus the value  $(r/m)^2$  for that particular replication. The reason we plot the MSE against  $(r/m)^2$  is that Theorem 1 predicts that the MSE of the slice learning algorithm should scale linearly with this value, which appears to be the case in Figure 4. On the other hand, the best theoretical results for the convex algorithm state that the MSE scales linearly with  $r/m$  (i.e., sublinear in  $(r/m)^2$ ); Figure 4 confirms that this is indeed the case. That is to say, the slice learning algorithm outperforms the convex algorithm in recovering tensors of low Tucker rank, even though the convex algorithm is *suites specifically for tensors of low Tucker rank*.

In terms of computation, the slice learning algorithm is also superior. The first-order Douglas–Rachford splitting method for solving (2) requires three SVDs in each iteration. Compared with the two SVDs required by the entire slice learning algorithm, this means *each iteration is slower than the entire slice learning algorithm*, and in our experiments, the whole algorithm was consistently at least 10 times slower. Along these lines, there is ongoing progress in improving the computational efficiency of convex optimization approaches, for example, by Liu et al. (2014).

**Figure 4.** (Color online) Comparison of Slice Learning and Convex Approach for Noisy Recovery of Low Tucker Rank Tensors



Note. MSE vs.  $(r/m)^2$  is plotted for each replication.

**6.1.2. The Value of Side Information.** We performed a similar experiment to test the recovery of tensors—particularly, recovery in terms of SMSE—with varying numbers of slices. We randomly generated  $m \times m \times n$  tensors of slice rank  $r$  where in each replication,  $m$  was drawn randomly between 10 and 50; then  $r$  was drawn randomly between 1 and  $m$ ; and finally  $n$  was set equal to either 1,  $r$ , or  $m$ . In each replication, we randomly drew orthonormal  $m \times r$  matrices  $U$  and  $V$ , along with  $r \times r$  matrices  $S^1, \dots, S^n$ , with entries drawn from the standard normal distribution, and we set the slices of the ground truth tensor to be  $M^k = US^kV^\top$ . Independent mean-zero Gaussian noise of standard deviation 0.1 was then added as before.

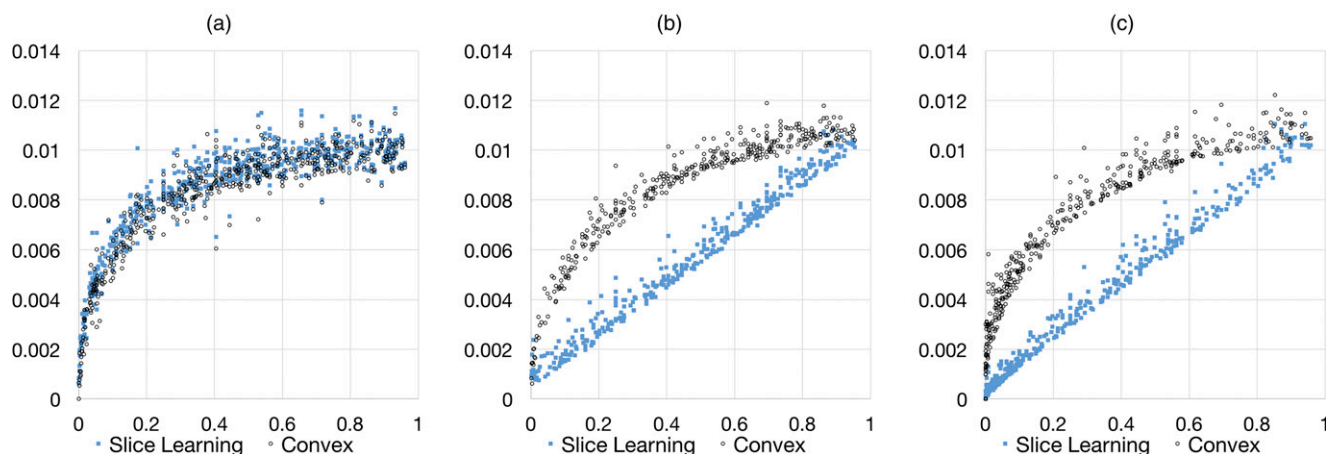
We used a similar convex algorithm as a benchmark:

$$\arg \min_Y \sum_{i=1}^2 \lambda \|Y_{(i)}\|_* + \|Y - X\|_F^2. \tag{6}$$

This objective is almost identical to (2), except that it does not include the nuclear norm of the mode 3 unfolding. This is catered to recovering low slice rank tensors, as it imposes no penalty on the complexity between slices. We solved this with a similar Douglas–Rachford splitting method that requires two singular value decompositions in each iteration. Just as in (2), the parameter  $\lambda$  encodes knowledge of rank, and so we solve (6) for  $\lambda$  ranging from  $2^{-2}$  to  $2^5$  and report the best performance among all of these.

We performed 400 replications. In each replication, we measured the SMSE of the slice learning algorithm and the convex algorithm (6). The results for each of the three cases ( $n = 1, r, m$ ) are depicted in three separate plots in Figure 5. Figure 5(a) shows that with a single slice, both algorithms have SMSE sublinear in  $(r/m)^2$ ; this is to be expected as the exercise is equivalent to matrix recovery where the best achievable rate is  $r/m$  (see Proposition 2). Figure 5(c) shows that the

**Figure 5.** (Color online) Comparison of Slice Learning and the Convex Approach for Noisy Slice Recovery of Low Slice Rank Tensors with Varying Numbers of Slices



Notes. SMSE vs.  $(r/m)^2$  is plotted for each replication. Panel (a):  $n = 1$ ; panel (b):  $n = r$ ; panel (c):  $n = m$ .

slice learning algorithm achieves the gold standard performance of SMSE linear in  $(r/m)^2$  with  $n = m$  slices, whereas the convex algorithm is still sublinear. The good performance of the slice learning algorithm is to be expected because in Theorem 1 we were able to show that  $n = m$  slices are sufficient to achieve the  $(r/m)^2$  rate. The surprising result is that Figure 5(b) is almost identical to Figure 5(c), implying that  $n = r$  slices are sufficient to achieve this same rate. This suggests that there are settings where slice learning can greatly outperform standard matrix learning, with only *very little* side information.

## 6.2. Experiments on Real Data

In addition to synthetic experiments, we also performed experiments using real-world data to address a number of important challenges that occur in practice. In particular, these experiments differ greatly from the previous synthetic experiments in that the data are very large and sparse. The data are from Xiami.com, a major online music streaming service where users may listen to songs and share their own music. Within the service, users can interact with songs in different ways: they can “listen” to, “download,” and “collect” any song offered by the service. The collect interaction is especially important, as it is a strong signal of a user’s affinity for a song, but it is performed with the least frequency in our data. Our data set<sup>4</sup> is a sample of all three interaction types between users and songs over a six-month period in 2015. For the experiments, we represent these data as a three-dimensional tensor  $X$  with three slices, one for each type of interaction, with binary entries indicating whether that particular user–song interaction occurred during the six-month period. Just as in our motivating retail example, we model this as Bernoulli noise; that is, we assume that the data  $X$  are a Bernoulli

observation of some ground truth tensor  $M$  of probabilities:  $X_{ij}^k \sim \text{Ber}(M_{ij}^k)$ .

We compared the slice learning algorithm against two benchmarks. The first benchmark is the *naïve* approach of using a matrix recovery algorithm separately on each slice, which is still a typical approach to collaborative filtering. The second benchmark is a more sophisticated approach: form one of the unfoldings and use a matrix recovery algorithm on the unfolding; we will refer to this algorithm as the *matrix* approach. This approach nicely exploits tensor structure and can be performed at large scale; Mu et al. (2013) study such algorithms and demonstrate theoretical guarantees for high-dimensional tensors.

Both benchmarks require a matrix recovery algorithm, and for our experiments, we recovered matrices by replacing hidden entries with zeros and then calculating a low-rank approximation of this modified matrix (Achlioptas and McSherry 2007, Keshavan et al. 2010). This procedure requires roughly the same computational budget as the slice learning algorithm. Finally, because our experiments were quite large (see Table 1), calculating SVDs using standard libraries meant for dense matrices was not feasible. Instead, we used the off-the-shelf software package PROPACK<sup>5</sup> (Larsen 1998), which exploits data sparsity through the iterative algorithms described in Section 4.1.

**6.2.1. Experimental Setup.** Because the data are binary, we evaluated performance vis-à-vis a binary classification task (i.e., the task of classifying entries as being equal to 0 or 1) on half of the entries. Because each entry of the tensor corresponds to whether a particular user interacted with a particular song, we will refer to the values 0 and 1 as “not occurring” and “occurring,” respectively.

**Table 1.** Summary of Experiments on Xiami Data for Recovering the Collect Slice

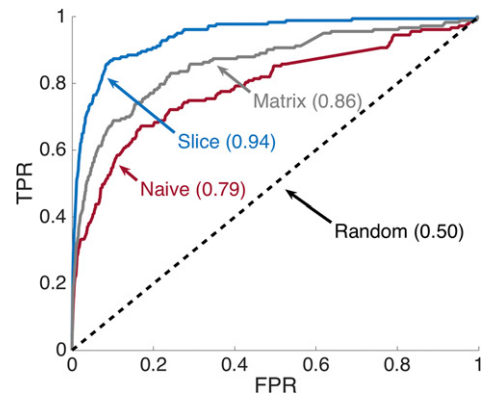
Users	Songs	Sparse	Naïve	Matrix	Slice
2,412	1,541	5.7	0.76 (4)	0.83 (7)	0.91 (14)
4,951	2,049	4.1	0.73 (7)	0.78 (12)	0.91 (15)
27,411	3,472	2.0	0.66 (9)	0.67 (19)	0.87 (20)
23,300	10,106	1.0	0.86 (1)	0.87 (1)	0.95 (18)
53,713	10,199	0.6	0.82 (3)	0.84 (1)	0.95 (13)

*Notes.* Each row corresponds to an experiment on a subset of the data. Columns “Users” and “Songs” show the number of users and songs in each experiment, respectively, and “Sparse” gives the average number of collects per user in the data. Results for the naïve benchmark, the matrix-based benchmark, and the slice learning algorithm are shown in the last three columns. The average AUC over 10 replications is reported, along with the rank in parentheses.

Our algorithm and the two benchmark algorithms return complete tensors with continuous values. To convert these continuous values to classifications of occurring and not occurring, we chose a fixed threshold  $\theta$  and classified any entry exceeding  $\theta$  as occurring and the remaining entries as not occurring. To evaluate an algorithm’s performance, we can calculate two important performance metrics: the true-positive rate (TPR) and the false-positive rate (FPR). Of all the hidden entries that occurred, the TPR is the proportion that the algorithm correctly classified as occurring, and of all the hidden entries that did not occur, the FPR is the proportion that the algorithm incorrectly classified as occurring.

A good classification has a high TPR and low FPR. This might inform the choice of the threshold  $\theta$ , but unfortunately, there is a trade-off: both the TPR and FPR are nonincreasing in  $\theta$ , so it is impossible to improve on both metrics just by changing  $\theta$ . In particular, by varying  $\theta$ , a given algorithm produces a set of values ranging from all entries being classified as not occurring (TPR and FPR are equal to 0) to all entries being classified as occurring (TPR and FPR are equal to 1). To evaluate the performance of an algorithm independent from the choice of  $\theta$ , we can plot each of these classifications in the receiver operating characteristic (ROC) space and calculate the area under the resulting curve (AUC), as in Figure 6. The AUC is always in  $[0, 1]$ , and a higher value generally signifies greater accuracy; as a benchmark, the AUC of a random classification is expected to be 0.5.

**6.2.2. Summary of Results.** The results of the experiment, in terms of recovering the collect slice, are summarized in Table 1. We performed experiments on tensors of five different sizes, described in the first two columns. These tensors were created by selecting subsets of the densest rows and columns of the original data set. The third column shows the sparsity of the collect slice in each of the five tensors, measured in the average number of collects per user. Note that the data

**Figure 6.** (Color online) Sample ROC Curves for Recovering the Collect Slice

*Notes.* These curves were generated from a single replication of the experiment in the first row of Table 1. For each algorithm, TPR vs. FPR is plotted, and the AUC is reported.

are extremely sparse—in the largest tensor, we observe less than a single collect per user.

For each tensor, 10 replications were performed, where each replication included a resampling of observed entries, followed by performances of all three algorithms with ranks ranging from 1 to 20. For each algorithm and rank, the AUC was averaged over all 10 replications. The best average AUC of each algorithm is reported in the last three columns of Table 1, and the best ranks are given in parentheses.

In absolute terms, the slice learning algorithm performs very well, with an AUC consistently above 0.87. The algorithm also consistently outperforms both benchmark approaches by a significant margin. For example, in the largest experiment with approximately 50K users and 10K songs, the slice learning algorithm has an average AUC of 0.95, whereas the naïve and matrix algorithms have AUCs of 0.82 and 0.84, respectively. Part of this strong performance comes from the fact that the slice learning algorithm is able to estimate more complex models, which is demonstrated by the consistently higher rank values.

The equivalent tables for recovering the listen and download slices can be found in Section F in the online appendix. To summarize those results succinctly, the slice learning algorithm again consistently outperforms both benchmarks across all experiments. The margin of improvement is lower, but that is to be expected, as those slices are less sparse than the collect slice, and so the collect slice does not offer as much side information.

## 7. Conclusion

This paper introduced a new approach to modeling and learning with side information. Motivated by settings in e-commerce, where data are sparse but multiple interactions occur, we formulated the problem of recovering slices of a three-dimensional tensor from a noisy



observation of the tensor. We proposed the slice learning algorithm, a computationally efficient algorithm that scales to enormous size by leveraging sparsity. From a theoretical standpoint, we showed that the algorithm achieves the minimax lower bound for recovery with sufficiently many slices; this guarantee is the best-known guarantee for noisy recovery of tensors and the first guarantee of its kind for recovery of specific slices. Synthetic experiments further supported the fact that this algorithm outperforms existing convex methods in both efficiency and accuracy. Experiments on real-world data from the music streaming service Xiami.com demonstrated the scalability of the approach and provided real empirical evidence that having side information is advantageous and that our approach utilizes side information effectively.

Our work points to a number of interesting, exciting directions for future work.

1. *Different forms of side information*: Side information may come in the form of data specific to the row space or the column space. For example, retailers have demographic information on their customers and basic information about their products. As we discussed in Section 4, the slice learning algorithm can incorporate this kind of side information. A deeper analysis of this procedure is an important next step.

2. *Trimming*: Our algorithm performs best under the balanced noise assumption. We defined precisely how to measure the level of unbalance and quantified the penalty of unbalance. When the noise is known to be significantly unbalanced, it may be possible to weight the rows and columns of the tensor in such a way that induces balanced noise. This weighted tensor can be estimated and then unweighted to recover the original tensor. Such “trimming” procedures need to be analyzed.

3. *Higher-dimensional tensors*: There are many applications for tensors of dimension greater than 3. For example, retailers might view their sales transactions over time, producing a three-dimensional tensor where time is the third dimension. Time-series data for multiple interactions may then be viewed as a four-dimensional tensor. There may be ways that the slice learning algorithm can be generalized to higher dimensions.

## Endnotes

<sup>1</sup> For completeness, we can also define a third unfolding: for each row and column index, we can take the corresponding entries across all the slices to create a column vector in  $\mathbb{R}^n$ ; these are the columns of the mode 3 unfolding, denoted as  $M_{(3)}$ , an  $n \times m^2$  matrix.

<sup>2</sup> In our analysis, we will place substantially weaker requirements on  $\epsilon$ .

<sup>3</sup> The condition number of a matrix can be defined as the ratio of its largest singular value to its smallest singular value. The ratio we describe here is defined only for the *nonzero* singular values.

<sup>4</sup> See <https://tianchi.shuju.aliyun.com/>, accessed December 20, 2016.

<sup>5</sup> See <http://sun.stanford.edu/~rmunk/PROPACK/>, accessed December 20, 2016.

## References

- Acar E, Çamtepe SA, Krishnamoorthy MS, Yener B (2005) Modeling and multiway analysis of chatroom tensors. Kantor P, Muresan G, Roberts F, Zeng D, Wang F-Y, Chen H, Merkle R, eds. *IEEE Internat. Conf. Intelligence Security Informatics Proc.* (Springer, Berlin), 256–268.
- Achlioptas D, McSherry F (2007) Fast computation of low-rank matrix approximations. *J. ACM* 54(2):Article 9.
- Adomavicius G, Tuzhilin A (2005) Toward the next generation of recommender systems: A survey of the state-of-the-art and possible extensions. *IEEE Trans. Knowledge Data Engrg.* 17(6):734–749.
- Ansari A, Mela CF (2003) E-customization. *J. Marketing Res.* 40(2): 131–145.
- Ansari A, Essegai S, Kohli R (2000) Internet recommendation systems. *J. Marketing Res.* 37(3):363–375.
- Bahadori MT, Yu QR, Liu Y (2014) Fast multivariate spatio-temporal analysis via low rank tensor learning. Ghahramani Z, Welling M, Cortes C, Lawrence ND, Weinberger KQ, eds. *Advances in Neural Information Processing Systems*, vol. 27 (Curran Associates, Red Hook, NY), 3491–3499.
- Bernhardsson E (2014) Music discovery at Spotify. Presentation, Machine Learning Conference, April 14, Helen Mills Event Center, New York.
- Berry MW, Dumais ST, O’Brien GW (1995) Using linear algebra for intelligent information retrieval. *SIAM Rev.* 37(4):573–595.
- Besbes O, Gur Y, Zeevi A (2015) Optimization in online content recommendation services: Beyond click-through rates. *Manufacturing Service Oper. Management* 18(1):15–33.
- Bodapati AV (2008) Recommendation systems with purchase data. *J. Marketing Res.* 45(1):77–93.
- Candès EJ, Plan Y (2011) Tight oracle inequalities for low-rank matrix recovery from a minimal number of noisy random measurements. *IEEE Trans. Inform. Theory* 57(4):2342–2359.
- Candès EJ, Tao T (2010) The power of convex relaxation: Near-optimal matrix completion. *IEEE Trans. Inform. Theory* 56(5):2053–2080.
- Carroll JD, Chang JJ (1970) Analysis of individual differences in multidimensional scaling via an n-way generalization of “Eckart-Young” decomposition. *Psychometrika* 35(3):283–319.
- Cattell RB (1966) The scree test for the number of factors. *Multivariate Behav. Res.* 1(2):245–276.
- Chatterjee S (2014) Matrix estimation by universal singular value thresholding. *Ann. Statist.* 43(1):177–214.
- Chiang K-Y, Hsieh C-J, Dhillon IS (2015) Matrix completion with noisy side information. Cortes C, Lawrence ND, Lee DD, Sugiyama M, Garnett R, eds. *Advances in Neural Information Processing Systems*, vol. 28 (Curran Associates, Red Hook, NY), 3429–3437.
- Chung TS, Rust RT, Wedel M (2009) My mobile music: An adaptive personalization system for digital audio players. *Marketing Sci.* 28(1):52–68.
- Demirezen EM, Kumar S (2016) Optimization of recommender systems based on inventory. *Production Oper. Management* 25(4):593–608.
- Fleder D, Hosanagar K (2009) Blockbuster culture’s next rise or fall: The impact of recommender systems on sales diversity. *Management Sci.* 55(5):697–712.
- Gandy S, Recht B, Yamada I (2011) Tensor completion and low-rank tensor recovery via convex optimization. *Inverse Problems* 27(2):Article 025010.
- Garfinkel R, Gopal R, Pathak B, Yin F (2008) Shopbot 2.0: Integrating recommendations and promotions with comparison shopping. *Decision Support Systems* 46(1):61–69.
- Gavish M, Donoho DL (2014) The optimal hard threshold for singular values is  $4/\sqrt{3}$ . *IEEE Trans. Inform. Theory* 60(8):5040–5053.
- Ghose A, Ipeirotis PG, Li B (2012) Designing ranking systems for hotels on travel search engines by mining user-generated and crowdsourced content. *Marketing Sci.* 31(3):493–520.

- Goldberg D, Nichols D, Oki BM, Terry D (1992) Using collaborative filtering to weave an information tapestry. *Comm. ACM* 35(12):61–70.
- Grey P (2015) How many products does Amazon sell? *Export-x* (December 11), <https://export-x.com/2015/12/11/how-many-products-does-amazon-sell-2015/>.
- Gross D (2011) Recovering low-rank matrices from few coefficients in any basis. *IEEE Trans. Inform. Theory* 57(3):1548–1566.
- Grover R, Srinivasan V (1987) A simultaneous approach to market segmentation and market structuring. *J. Marketing Res.* 24(2): 139–153.
- Harris A (2015) Amazon sold 5 billion items in 2014. *Fast Company* (January 6), <http://www.fastcompany.com/3040445/fast-feed/amazon-sold-5-billion-items-in-2014>.
- Henrion R (1994) N-way principal component analysis theory, algorithms and applications. *Chemometrics Intelligent Laboratory Systems* 25(1):1–23.
- Hoff PD (2007) Model averaging and dimension selection for the singular value decomposition. *J. Amer. Statist. Assoc.* 102(478): 674–685.
- Hosanagar K, Krishnan R, Ma L (2008) Recommended for you: The impact of profit incentives on the relevance of online recommendations. *ICIS 2008 Proc.* (Association for Information Systems, Atlanta), Paper 31.
- Huang B, Mu C, Goldfarb D, Wright J (2015) Provable models for robust low-rank tensor completion. *Pacific J. Optim.* 11(2): 339–364.
- Jacobs BJ, Donkers B, Fok D (2016) Model-based purchase predictions for large assortments. *Marketing Sci.* 35(3):389–404.
- Jain P, Dhillon IS (2013) Provable inductive matrix completion. Working paper, University of Texas at Austin, Austin.
- Jain P, Oh S (2014) Provable tensor factorization with missing data. Ghahramani Z, Welling M, Cortes C, Lawrence ND, Weinberger KQ, eds. *Advances in Neural Information Processing Systems*, vol. 27 (Curran Associates, Red Hook, NY), 1431–1439.
- Keshavan RH, Montanari A, Oh S (2010) Matrix completion from a few entries. *IEEE Trans. Inform. Theory* 56(6):2980–2998.
- Kleinberg JM (1999) Authoritative sources in a hyperlinked environment. *J. ACM* 46(5):604–632.
- Kolda TG, Bader BW (2009) Tensor decompositions and applications. *SIAM Rev.* 51(3):455–500.
- Koltchinskii V, Lounici K, Tsybakov AB (2011) Nuclear-norm penalization and optimal rates for noisy low-rank matrix completion. *Ann. Statist.* 39(5):2302–2329.
- Koren Y, Bell R, Volinsky C (2009) Matrix factorization techniques for recommender systems. *Computer* 42(8):30–37.
- Larsen RM (1998) Lanczos bidiagonalization with partial reorthogonalization. *DAIMI Rep. Ser.* 27(537).
- Liben-Nowell D, Kleinberg J (2007) The link-prediction problem for social networks. *J. Amer. Soc. Inform. Sci. Techn.* 58(7):1019–1031.
- Linden G, Smith B, York J (2003) Amazon.com recommendations: Item-to-item collaborative filtering. *IEEE Internet Comput.* 7(1): 76–80.
- Liu J, Musialski P, Wonka P, Ye J (2013) Tensor completion for estimating missing values in visual data. *IEEE Trans. Pattern Anal. Machine Intelligence* 35(1):208–220.
- Liu Y, Shang F, Fan W, Cheng J, Cheng H (2014) Generalized higher-order orthogonal iteration for tensor decomposition and completion. Ghahramani Z, Welling M, Cortes C, Lawrence ND, Weinberger KQ, eds. *Advances in Neural Information Processing Systems*, vol. 27 (Curran Associates, Red Hook, NY), 1763–1771.
- Marčenko VA, Pastur LA (1967) Distribution of eigenvalues for some sets of random matrices. *Math. USSR-Sbornik* 1(4):457–483.
- Mørup M (2011) Applications of tensor (multiway array) factorizations and decompositions in data mining. *Wiley Interdisciplinary Rev.: Data Mining Knowledge Discovery* 1(1):24–40.
- Mu C, Huang B, Wright J, Goldfarb D (2013) Square deal: Lower bounds and improved relaxations for tensor recovery. Working paper, Columbia University, New York.
- Naik P, Wedel M, Bacon L, Bodapati A, Bradlow E, Kamakura W, Kreulen J, Lenk P, Madigan DM, Montgomery A (2008) Challenges and opportunities in high-dimensional choice data analyses. *Marketing Lett.* 19(3–4):201–213.
- Oliver NM, Rosario B, Pentland AP (2000) A Bayesian computer vision system for modeling human interactions. *IEEE Trans. Pattern Anal. Machine Intelligence* 22(8):831–843.
- Owen AB, Perry PO (2009) Bi-cross-validation of the SVD and the nonnegative matrix factorization. *Ann. Appl. Statist.* 3(2): 564–594.
- Papadimitriou CH, Tamaki H, Raghavan P, Vempala S (1998) Latent semantic indexing: A probabilistic analysis. *Proc. 17th ACM SIGACT-SIGMOD-SIGART Sympos. Principles Database Systems* (ACM, New York), 159–168.
- Recht B (2011) A simpler approach to matrix completion. *J. Machine Learn. Res.* 12(February):3413–3430.
- Romera-Paredes B, Pontil M (2013) A new convex relaxation for tensor completion. Burges CLC, Bottou L, Welling M, Ghahramani Z, Weinberger KQ, eds. *Advances in Neural Information Processing Systems*, vol. 26 (Curran Associates, Red Hook, NY), 2967–2975.
- Romera-Paredes B, Aung H, Bianchi-Berthouze N, Pontil M (2013) Multilinear multitask learning. *Proc. Machine Learn. Res.* 28(3): 1444–1452.
- Shi Y, Larson M, Hanjalic A (2014) Collaborative filtering beyond the user-item matrix: A survey of the state of the art and future challenges. *ACM Comput. Surveys* 47(1):Article 3.
- Signoretto M, Dinh QT, De Lathauwer L, Suykens JA (2014) Learning with tensors: A framework based on convex optimization and spectral regularization. *Machine Learn.* 94(3):303–351.
- Singh AP, Gordon GJ (2008) Relational learning via collective matrix factorization. *Proc. 14th ACM SIGKDD Internat. Conf. Knowledge Discovery Data Mining* (ACM, New York), 650–658.
- Sirovich L, Kirby M (1987) Low-dimensional procedure for the characterization of human faces. *J. Optical Soc. Amer. A* 4(3): 519–524.
- So AMC, Ye Y (2007) Theory of semidefinite programming for sensor network localization. *Math. Programming* 109(2–3):367–384.
- Soni A, Jain S, Haupt J, Gonella S (2016) Noisy matrix completion under sparse factor models. *IEEE Trans. Inform. Theory* 62(6): 3636–3661.
- Statista (2016) Annual number of worldwide active Amazon customer accounts from 1997 to 2015. Accessed December 20, 2016, <http://www.statista.com/statistics/237810/number-of-active-amazon-customer-accounts-worldwide/>.
- Sun JT, Zeng HJ, Liu H, Lu Y, Chen Z (2005) CubeSVD: A novel approach to personalized web search. *Proc. 14th Internat. Conf. World Wide Web* (ACM, New York), 382–390.
- Suzuki T (2015) Convergence rate of Bayesian tensor estimator and its minimax optimality. *Proc. Machine Learn. Res.* 37:1273–1282.
- Tomioka R, Suzuki T (2013) Convex tensor decomposition via structured Schatten norm regularization. Burges CLC, Bottou L, Welling M, Ghahramani Z, Weinberger KQ, eds. *Advances in Neural Information Processing Systems*, vol. 26 (Curran Associates, Red Hook, NY), 1331–1339.
- Tomioka R, Suzuki T, Hayashi K, Kashima H (2011) Statistical performance of convex tensor decomposition. Shawe-Taylor J, Zemel RS, Bartlett PL, Pereira F, Weinberger KQ, eds. *Advances in Neural Information Processing Systems*, vol. 24 (Curran Associates, Red Hook, NY), 972–980.
- Troyanskaya O, Cantor M, Sherlock G, Brown P, Hastie T, Tibshirani R, Botstein D, Altman RB (2001) Missing value estimation methods for DNA microarrays. *Bioinformatics* 17(6):520–525.

- Vasilescu MAO, Terzopoulos D (2002) Multilinear analysis of image ensembles: TensorFaces. Heyden A, Sparr G, Nielsen M, Johansen P, eds. *Computer Vision—ECCV 2002*, Lecture Notes in Computer Science, vol. 2350 (Springer, Berlin), 447–460.
- Wold S (1978) Cross-validated estimation of the number of components in factor and principal components models. *Technometrics* 20(4):397–405.
- Xu M, Jin R, Zhou ZH (2013) Speedup matrix completion with side information: Application to multi-label learning. Burges CLC, Bottou L, Welling M, Ghahramani Z, Weinberger KQ, eds. *Advances in Neural Information Processing Systems*, vol. 26 (Curran Associates, Red Hook, NY), 2301–2309.
- Yin YQ, Bai ZD, Krishnaiah PR (1988) On the limit of the largest eigenvalue of the large dimensional sample covariance matrix. *Probab. Theory Related Fields* 78(4):509–521.
- Yuan M, Zhang CH (2015) On tensor completion via nuclear norm minimization. *Foundations Comput. Math.* 16(4):1031–1068.
- Zheng Q, Tomioka R (2015) Interpolating convex and non-convex tensor decompositions via the subspace norm. Cortes C, Lawrence ND, Lee DD, Sugiyama M, Garnett R, eds. *Advances in Neural Information Processing Systems*, vol. 28 (Curran Associates, Red Hook, NY), 3106–3113.
- Zhou H, Li L, Zhu H (2013) Tensor regression with applications in neuroimaging data analysis. *J. Amer. Statist. Assoc.* 108(502): 540–552.