

# Manipulation Robustness of Collaborative Filtering Systems

Xiang Yan  
Stanford University  
xyan@stanford.edu

Benjamin Van Roy  
Stanford University  
bvr@stanford.edu

April 30, 2009

## Abstract

A collaborative filtering system recommends to users products that similar users like. Collaborative filtering systems influence purchase decisions, and hence have become targets of manipulation by unscrupulous vendors. We provide theoretical and empirical results demonstrating that while common nearest neighbor algorithms, which are widely used in commercial systems, can be highly susceptible to manipulation, two classes of collaborative filtering algorithms which we refer to as *linear* and *asymptotically linear* are relatively robust. These results provide guidance for the design of future collaborative filtering systems.

## 1 Introduction

While the expanding universe of products available via Internet commerce provides consumers with valuable options, sifting through the numerous alternatives to identify desirable choices can be challenging. Collaborative filtering (CF) systems aid this process by recommending to users products desired by similar individuals.

At the heart of a CF system is an algorithm that predicts whether a given user will like various products based on his past behavior and that of other users. Nearest neighbor (NN) algorithms, for example, have enjoyed wide use in commercial CF systems, including those of Amazon, Netflix, and Youtube (Bennett, 2006; Linden et al., 2003; Ryan, 2008). A prototypical NN algorithm stores each user's history, which may include, for instance, his product ratings and purchase decisions. To predict whether a particular user will like a particular product, the algorithm identifies a number of other users with similar histories. A prediction is then generated based on how these so-called neighbors have responded to the product. This prediction could be, for example, a weighted average of past ratings supplied by neighbors.

Because purchase decisions are influenced by CF systems, they have become targets of manipulation by unscrupulous vendors. For instance, a vendor can create multiple online identities and use each to rate his own product highly and competitors' products poorly. As an example, Amazon's CF system was manipulated so that users who viewed a spiritual guide written by a well-known Christian evangelist were subsequently recommended a sex manual for gay men (Olsen, 2002). Although this incident may not have been driven by commercial motives, it highlights the vulnerability of CF systems. The research literature offers further empirical evidence that NN algorithms are susceptible to manipulation (Burke et al., 2005; Lam and Riedl, 2004; Mehta and Nejdl, 2008; Mobasher et al., 2005, 2006; O'Mahony et al., 2004; Sandvig et al., 2007; Zhang et al., 2006).

In order to curb manipulation, one might consider authenticating each user by asking for, say, a credit card number to limit the number of fake identities. This may be effective in some situations. However, in web services that do not facilitate financial transactions, such as Youtube, requiring authentication would intrude privacy and drive users away. One might also consider using only customer purchase data, when they are available, as a basis for recommendations because they are likely generated by honest users. Recommendation quality may be improved, however, if higher-volume data such as page views are also properly utilized.

In this paper, we seek to understand the extent to which manipulators can hurt the performance of CF systems and how CF algorithms should be designed to abate their influence. We find that, while NN algorithms can be quite sensitive to manipulation, CF algorithms that carry out predictions based on a particular class of probabilistic models are surprisingly robust. For reasons that we will explain in the paper, we will refer to algorithms of this kind as *linear CF algorithms*.

We find that as a user rates an increasing number of products, the average accuracy of predictions made by a linear CF algorithm becomes insensitive to manipulated data. For instance, even if half of all ratings are provided by manipulators who try to promote half of the products, predictions for users with long histories will barely be distorted, on average. To provide some intuition for why our results should hold, we now offer an informal argument. A robust CF algorithm should learn from its mistakes. In particular, differences between its predictions and actual ratings should help improve predictions on future ratings. A linear CF algorithm generates predictions based on a

probability distribution that is a convex combination of two distributions: one that it would learn given only data generated by honest users and one that it would learn given only manipulated data. As a user whose ratings we wish to predict provides more ratings, it becomes increasingly clear which of these two distributions better represents his preferences. As a result, the weight placed on manipulated data diminishes and distortion vanishes.

The main theoretical result of this paper formalizes the above argument. In particular, we will define a notion of distortion induced by manipulators and establish an upper bound on distortion, which takes a particularly simple form:

$$\text{distortion} \leq \frac{1}{n} \ln \frac{1}{1-r}.$$

Here  $r$  is the fraction of data that is generated by manipulators and  $n$  is the number of products that have already been rated by a user whose future ratings we wish to predict. The bound is very general. First, it applies to all linear CF algorithms. Second, it applies to all manipulation strategies even if manipulators coordinate their actions and produce data with knowledge of all data generated by honest users. The bound demonstrates that as the number of prior ratings  $n$  increases, distortion vanishes. It also identifies the number required to limit distortion to a certain level. This offers guidance for the design of a recommendation system: the system may, for example, assess and inform users about the confidence of each recommendation. The system may also require a new user to rate a set number of products before making recommendations to him. To put this in perspective, consider the following numerical example. Suppose a CF system that accepts binary ratings predicts future ratings correctly 80% of the time in the absence of manipulation. If 10% of all ratings are provided by manipulators, according to our bound, the system can maintain a 75% rate of correct predictions by requiring each new user to rate at least 21 products before receiving recommendations.

To broaden the scope of our analysis, we will also study CF algorithms that behave like linear CF algorithms asymptotically as the size of the training set grows. This class of algorithms, which we refer to as *asymptotically linear*, is more flexible in accommodating modeling assumptions that may improve prediction accuracy. We will establish that a relaxed version of our distortion bound

for linear CF algorithms applies to asymptotically linear CF algorithms.

We will also show that our distortion bound does not generally hold for NN algorithms. Intuitively, this is because prediction errors do not always improve the selection of neighbors. In particular, as a user provides more ratings, manipulated data that contribute to inaccurate predictions of his future ratings may remain in the set of neighbors while data generated by honest users may be eliminated from it. As a result, distortion of predictions may not decrease. We will later provide an example to illustrate this.

In addition to theoretical results, this paper provides an empirical analysis using a publicly available set of movie ratings generated by users of Netflix’s recommendation system. We produce a distorted version of this data set by injecting manipulated ratings generated using a manipulation technique studied in prior literature. We then compare results from application of three CF algorithms: an NN algorithm, a linear CF algorithm called the kernel density estimation algorithm, and an asymptotically linear CF algorithm called the naive Bayes algorithm. Results demonstrate that while performance of the NN algorithm is highly susceptible to manipulation, those of kernel density estimation and naive Bayes algorithms are relatively robust. In particular, the latter two experience distortion lower than the theoretical bound we provide, whereas the distortion for the former exceeds it by far.

One might also wonder whether manipulation robustness of a CF algorithm comes at the expense of its prediction accuracy. As an example, consider an algorithm that fixes predictions for all ratings to be a constant, without regard to the training data. This algorithm is uninfluenced by manipulation but is likely to yield poor predictions, and is therefore not useful. In our experiments, the accuracy demonstrated by the three algorithms all seems reasonable. This suggests that accuracy of a CF algorithm may be achieved alongside robustness.

Our theoretical and empirical results together suggest that commercial recommendation systems using NN algorithms can be made more robust by adopting approaches that we describe. Note that we are not proposing that real-world systems should implement the specific algorithms we present in this paper. Rather, our analysis highlights properties of CF algorithms that lead to robustness and practitioners may benefit from taking these properties into consideration when designing CF systems.

This paper is organized as follows. In the next section, we discuss some related work. In Section 3, we formulate a simplified model that serves as a context for studying alternative CF algorithms. We then establish results concerning the manipulation robustness of NN, linear, and asymptotically linear CF algorithms in Section 4. In Section 5, we present our empirical study. We make some closing remarks in a final section.

## 2 Related Work

Early research on CF systems focused on their performance in the absence of manipulation (Breese et al., 1998; Drineas et al., 2002; Herlocker et al., 1999; Kleinberg and Sandler, 2008; Moon and Russell, 2008; Motwani and Vassilvitskii, 2007; Sarwar et al., 2001; Schafer et al., 2001). Almost all work on manipulation robustness has been empirical. For example, Burke et al. (2005); Lam and Riedl (2004); Mehta and NejdI (2008); Mobasher et al. (2005, 2006); O’Mahony et al. (2004); Sandvig et al. (2007); Zhang et al. (2006) present studies on product ratings made publicly available by Internet commerce sites. In each case, manipulated ratings were injected, and CF algorithms were tested on the altered data sets. The results point out that NN algorithms and their variants are susceptible to manipulation. This line of work identifies an effective manipulation scheme, which is to create multiple identities and with each identity, provide positive ratings on products to be promoted while rating other products in a manner indistinguishable from that of honest users. In Mehta and NejdI (2008); Mobasher et al. (2006); Zhang et al. (2006), algorithms based on probabilistic latent semantic analysis and principal component analysis were tested. It turns out that these algorithms are asymptotically linear under certain assumptions about the data, and indeed, empirical results in these papers suggest that they are relatively robust to manipulation. These prior results support the conclusions of our work.

To the best of our knowledge, the only prior theoretical work on manipulation robustness of CF algorithms is reported in O’Mahony et al. (2004). This work analyzed an NN algorithm that uses the majority rating among a set of neighbors as the prediction of a user’s rating in an asymptotic regime of many users, each of whom rates all products. Manipulators rate as honest users would except on one fixed product. A bound is established on the algorithm’s prediction error for this

product’s rating as a function of the percentage of ratings provided by manipulators. In our work, we do not require users to rate all products and do not constrain manipulators to any particular strategies. Further, we study the performance distortion on average, rather than for a single product. Finally, a primary contribution of our work is in establishing manipulation robustness of linear and asymptotically linear CF algorithms, which turns out to be superior to NN algorithms in this dimension.

Several researchers have proposed alternative approaches to abating the influence of manipulators. In Resnick and Sami (2007), a mechanism is proposed where users accumulate reputations while providing ratings that are later validated by observed product quality, and a user’s influence on ratings predictions is limited by his reputation. In this mechanism, a bound is established on the distortion induced by any finite number of manipulators. In Massa and Avesani (2008); O’Donovan and Smyth (2006), researchers propose leveraging trust relationships among users to weight recommendations and fend off manipulation. Mehta (2007); Mobasher et al. (2007); Sandvig et al. (2007); Williams et al. (2007) suggest detecting manipulated ratings based on their patterns and discounting their impact. Our work complements this growing literature. First, additional sources of information can be integrated into the probabilistic framework that we introduce in this paper to further enhance manipulation robustness. Second, the analytical methods that we develop may be useful for studying the benefits of incorporating such information.

Distortion due to manipulation may also be viewed as a loss of utility in a sequential decision problem induced by errors in initial beliefs. Our analysis is based on ideas similar to those that have been used to study the latter topic, which is discussed in Gossner and Tomala (2008).

More broadly speaking, apart from collaborative filtering, there are other ways to aggregate users’ response to products in order to provide recommendations. Research has been performed on the manipulation robustness of these systems as well. To get a flavor of this line of work, see Bhattacharjee and Goel (2007); Dellarocas (2006); Friedman et al. (2007); Miller et al. (2005).

### 3 Model

We now formulate a simplified model that will serve as a context for assessing performance of alternative CF algorithms. We will first define the product ratings that we work with and then introduce measures of distortion induced by manipulators. For the convenience of the reader, we summarize our mathematical notation in a table in Appendix B.

#### 3.1 Ratings Vectors

In our model, a user selects ratings from a set  $\bar{S}$ . To simplify our discussion, we let  $\bar{S}$  be a finite subset of  $[0, 1]$ . For example,  $\bar{S}$  could be  $\{0, 1\}$  with 0 representing a negative rating and 1 representing a positive rating. Note that all the results in this paper can be easily generalized to accommodate any finite set  $\bar{S}$ . There are  $N$  products, and a user's type is identified by a vector in  $\bar{S}^N$ . Each  $n$ th component of this vector reflects how the user would rate the  $n$ th product after inspecting it.

The CF system has access to ratings provided by  $M$  identities who have rated products in the past. The data from each  $m$ th identity takes the form of a ratings vector  $w^m \in S^N$ , where  $S = \bar{S} \cup \{?\}$ . Here, an element of  $\bar{S}$  represents a product rating whereas a question mark indicates that a product has not been rated. We refer to  $W = (w^1, \dots, w^M) \in S^{N \times M}$  as the *training data*. This data is used by a CF algorithm to predict future ratings.

Consider a user who is distinct from identities that generated the training data and for whom we will generate recommendations. We will refer to such a user as an *active user*. We will think of a CF algorithm as providing a probability mass function (PMF)  $p_{n,x,W}$  over  $\bar{S}$  for each triplet  $(n, x, W) \in \{1, \dots, N\} \times S^N \times S^{N \times M}$ . The PMF  $p_{n,x,W}$  represents beliefs about how an active user who has so far provided ratings  $x$  would rate product  $n$  after inspecting it. Such an algorithm can be used to guide recommendations; for example, the CF system might recommend to the active user the product he is most likely to rate highly among those that he has not already rated.

### 3.2 Distortion Measures

To study the influence of manipulation, we consider a situation where a fraction  $r$  of the identities are created by manipulators, while the remaining fraction  $1 - r$  correspond to distinct honest users. We denote the honest ratings vectors by  $y^1, \dots, y^{(1-r)M} \in S^N$  and the manipulated ratings vectors by  $z^1, \dots, z^{rM} \in S^N$ . Let  $Y = (y^1, \dots, y^{(1-r)M})$  and  $Z = (z^1, \dots, z^{rM})$  so that the training data is  $W = (Y, Z)$ .

To assess distortion of predictions made by a CF algorithm, we consider the following thought experiment. A hypothetical active user begins with a ratings vector  $x^0$ , with each  $n$ th component set to  $x_n^0 = ?$ , and inspects products in an order  $\nu = (\nu_1, \dots, \nu_N) \in \sigma_N$ , where  $\sigma_N$  denotes the set of permutations of  $\{1, \dots, N\}$ . After inspecting product  $\nu_k$ , the user rates it by sampling from the PMF  $p_{\nu_k, x^{k-1}, Y}$ . An updated ratings vector  $x^k$  is generated by incorporating this new rating in  $x^{k-1}$ . This stochastic process reflects how we would think honest users behave based on the CF algorithm and uncorrupted data set  $Y$ . We introduce the following measure of distortion, which we refer to as *Kullback-Leibler (KL) distortion*:

$$d_n^{\text{KL}}(p, \nu, Y, Z) = \frac{1}{n} \sum_{k=1}^n \mathbb{E} \left[ D \left( p_{\nu_k, x^{k-1}, Y} \parallel p_{\nu_k, x^{k-1}, (Y, Z)} \right) \right],$$

where  $D$  denotes Kullback-Leibler divergence with the natural log. That is, for any two PMFs  $p$  and  $q$  over support  $U$ ,  $D(p \parallel q) = \sum_{u \in U} p(u) \ln(p(u)/q(u))$ . This measure of the difference between PMFs is commonly used in information theory.

For each  $k$ , the PMF  $p_{\nu_k, x^{k-1}, Y}$  represents the prediction that would be made in the absence of manipulators, whereas  $p_{\nu_k, x^{k-1}, (Y, Z)}$  is what it becomes as a consequence of manipulation. Hence,  $D \left( p_{\nu_k, x^{k-1}, Y} \parallel p_{\nu_k, x^{k-1}, (Y, Z)} \right)$  measures the extent to which the manipulated data  $Z$  influences the prediction. We take the expectation of this quantity, with  $x^{k-1}$  distributed as the CF algorithm would have predicted if the data set were not corrupted by manipulated data. KL distortion  $d_n^{\text{KL}}(p, \nu, Y, Z)$  averages these terms over the first  $n$  inspected products.

Some algorithms such as NN algorithms generate predictions not in the form of PMFs, but as scalars that may be interpreted as the means of PMFs. For these algorithms, it may be more

suitable to measure manipulation impact in terms of *root-mean-squared (RMS) distortion*:

$$d_n^{\text{RMS}}(p, \nu, Y, Z) = \sqrt{\frac{1}{n} \sum_{k=1}^n \mathbb{E} \left[ \left( \tilde{x}_{\nu_k, x^{k-1}, Y} - \tilde{x}_{\nu_k, x^{k-1}, (Y, Z)} \right)^2 \right]},$$

where  $\tilde{x}_{\nu_k, x^{k-1}, Y}$  and  $\tilde{x}_{\nu_k, x^{k-1}, (Y, Z)}$  denote the scalar predictions of  $\bar{x}_{\nu_k}$  by the algorithm based on ratings history  $x^{k-1}$  and data sets  $Y$  and  $(Y, Z)$ , respectively. Note that if the algorithm generates PMFs as predictions,  $\tilde{x}_{\nu_k, x^{k-1}, Y}$  and  $\tilde{x}_{\nu_k, x^{k-1}, (Y, Z)}$  would be expectations of  $\bar{x}_{\nu_k}$  taken with respect to  $p_{\nu_k, x^{k-1}, Y}$  and  $p_{\nu_k, x^{k-1}, (Y, Z)}$ , respectively. The expectation in the definition of RMS distortion is taken with  $x^{k-1}$  distributed as the CF algorithm would have predicted based on  $Y$ . RMS distortion may offer a more transparent assessment than KL distortion because the former computes how much scalar predictions change in the same unit as the predictions themselves. RMS distortion is bounded by a function of KL distortion:

$$d_n^{\text{RMS}}(p, \nu, Y, Z) \leq \sqrt{\frac{1}{2} d_n^{\text{KL}}(p, \nu, Y, Z)}.$$

This is shown in Proposition 1 in Appendix A.1.

To offer an intuitive interpretation for RMS distortion, we consider a setting where users provide binary ratings and the CF system offers binary predictions based on the PMFs that it generates. That is, we set  $\bar{S} = \{0, 1\}$ . Given training data  $Y$ , for a user with ratings history  $x^{k-1}$ , the system generates a prediction of  $\tilde{x}_{\nu_k, x^{k-1}, Y} = 1$  for product  $\nu_k$  if  $p_{\nu_k, x^{k-1}, Y}(1) \geq 1/2$  and generates a prediction of  $\tilde{x}_{\nu_k, x^{k-1}, Y} = 0$  otherwise. Similarly, we denote  $\tilde{x}_{\nu_k, x^{k-1}, (Y, Z)}$  as the binary prediction based on  $(Y, Z)$ . We define the following *binary prediction distortion*:

$$d_n^{\text{B}}(p, \nu, Y, Z) = \frac{1}{n} \sum_{k=1}^n \left( \Pr(\bar{x}_{\nu_k} = \tilde{x}_{\nu_k, x^{k-1}, Y}) - \Pr(\bar{x}_{\nu_k} = \tilde{x}_{\nu_k, x^{k-1}, (Y, Z)}) \right),$$

where each  $x_{\nu_k}$  is distributed according to  $p_{\nu_k, x^{k-1}, Y}$  and  $x^{k-1}$  is distributed as the CF algorithm would have predicted based on  $Y$ . This quantity captures the average decrease in the probability of correct predictions, induced by manipulation. It turns out that binary prediction distortion is

bounded by RMS distortion:

$$d_n^B(p, \nu, Y, Z) \leq d_n^{\text{RMS}}(p, \nu, Y, Z). \quad (1)$$

Proved in Proposition 2 in Appendix A.1, this result offers an interpretation of RMS distortion as an upper bound on the drop in the probability of correct predictions in a binary setting.

One might wonder why we choose our particular distortion measures over other candidates. For instance, one option is to consider the top  $n$  most desirable products based on predictions, and define as distortion some measure of their quality change due to the manipulated samples. One reason why we prefer KL and RMS distortions is that they are convex functions of predictions while this measure is not. As such, this measure is difficult to analyze. Further, as will be discussed in Section 5.6, in a recent competition of CF algorithms, Netflix uses RMS error to assess their prediction accuracies (Netflix Prize, 2006). This suggests that commercial CF algorithms are typically designed to minimize convex measures of error. Our choice of distortion measures is in line with this approach.

Another option one might consider is to measure the worst distortion over all products. While it would be attractive to keep this measure small for all manipulation schemes, that may not be achievable. Further, a large worst-case distortion does not necessarily imply that a CF algorithm is not useful. Consider, for instance, a case where many manipulators aim to distort ratings of many different products, and the output of the algorithm is such that the rating of one product is significantly distorted while the rest are not. We might still say that the algorithm is reasonably robust whereas a worst-case measure would indicate otherwise. That said, we note that since KL and RMS distortions characterize average distortions, the robustness results of this paper do not provide guarantees on the distortion of individual products' ratings.

## 4 Collaborative Filtering Algorithms

In this section, we first introduce the notion of *probabilistic CF algorithms*. We then describe two classes of such algorithms, namely linear and asymptotically linear CF algorithms, and analyze their robustness to manipulation. Finally, we discuss nearest neighbor algorithms and their susceptibility to manipulation.

## 4.1 Probabilistic Collaborative Filtering Algorithms

A probabilistic CF algorithm carries out predictions based on a probabilistic model of how the training data is generated. We will model training data as being generated in the following way. First, user types  $\bar{w}^m \in \bar{S}^N$  are sampled i.i.d. from some PMF. Then,  $w^m \in S^N$  is sampled from a conditional PMF, conditioned on  $\bar{w}^m$ , which for each  $n$  assigns either  $w_n^m = ?$  or  $w_n^m = \bar{w}_n^m$ . Note that this model allows for dependence between the type of a user and the products he chooses to rate. This accommodates, for example, systems in which users tend to inspect and rate only products that they care for. Given a PMF  $\psi$  over  $\bar{S}^N \times S^N$ , we denote by  $\psi_{\bar{S}}$  and  $\psi_S$  the marginal PMFs over  $\bar{S}^N$  and  $S^N$ , respectively.

We will call a CF algorithm  $p$  *probabilistic* if for each  $W$  there exists a PMF  $\hat{\psi}^{p,W}$  over  $\bar{S}^N \times S^N$  such that for each  $n$  and  $x$ ,  $p_{n,x,W}$  is the marginal PMF of  $\bar{x}_n$  conditioned on  $x$ , with respect to the joint PMF  $\hat{\psi}^{p,W}$ . From here on, we will denote by  $\hat{\psi}_{\bar{S}}^{p,W}$  the marginal PMF over  $\bar{S}^N$  of the joint PMF  $\hat{\psi}^{p,W}$  corresponding to a probabilistic CF algorithm  $p$  and training set  $W$ .

## 4.2 Linear Collaborative Filtering Algorithms

We say that a probabilistic CF algorithm  $p$  is *linear* if for any  $W_1 \in S^{N \times M_1}$  and  $W_2 \in S^{N \times M_2}$ ,

$$\hat{\psi}_{\bar{S}}^{p,(W_1,W_2)} = \frac{M_1}{M_1 + M_2} \hat{\psi}_{\bar{S}}^{p,W_1} + \frac{M_2}{M_1 + M_2} \hat{\psi}_{\bar{S}}^{p,W_2}.$$

This definition states that the PMF  $\hat{\psi}_{\bar{S}}^{p,(W_1,W_2)}$  that a linear CF algorithm  $p$  generates based on training data  $(W_1, W_2)$  is a convex combination of two PMFs: namely, the PMF  $\hat{\psi}_{\bar{S}}^{p,W_1}$  that it generates based on  $W_1$  and the PMF  $\hat{\psi}_{\bar{S}}^{p,W_2}$  that it generates based on  $W_2$ .

We now examine the KL distortion that manipulators can induce on a linear CF algorithm. Consider training data  $W = (Y, Z)$  consisting of ratings vectors  $Y$  from honest users and  $Z$  from manipulators, with the latter making up a fraction  $r$  of the training data. The following theorem, which is the main theoretical contribution of this paper, establishes a bound on the resulting KL distortion.

**Theorem 1.** *Fix the number of products  $N$  and let  $p$  be a linear CF algorithm. Then, for all  $M$ ,*

$r \in \{0, 1/M, \dots, (M-1)/M\}$ ,  $Y \in S^{N \times (1-r)M}$ ,  $Z \in S^{N \times rM}$ , and  $\nu \in \sigma_N$ ,

$$d_n^{\text{KL}}(p, \nu, Y, Z) \leq \frac{1}{n} \ln \frac{1}{1-r}.$$

This result is proved in Appendix A.2.

Note that the bound only depends on the number of active user ratings  $n$  and the fraction of data  $r$  generated by manipulators. Hence, it represents a worst case bound over all linear CF algorithms  $p$ , the number of products  $N$ , the quantity  $M$  and values  $(Y, Z)$  of the training data, and the order  $\nu$  in which the active user rates products. This means, for example, that it applies even if manipulators coordinate with each other and select ratings with knowledge of the specific CF algorithm  $p$ , the honest ratings  $Y$ , and the ordering  $\nu$ . This also makes the bound relevant for realistic models of how a recommendation system might sequence products for a user; for example, each  $\nu_k$  could be the product that the CF algorithm predicts as being most desirable among remaining ones after the user has inspected products  $\nu_1, \dots, \nu_{k-1}$ .

Note that KL distortion vanishes as the number  $n$  of products rated by the active user increases. To develop intuition for why this happens, we now offer an informal argument. Observe that  $\psi_{\bar{S}}^{p, (Y, Z)} = (1-r)\psi_{\bar{S}}^{p, Y} + r\psi_{\bar{S}}^{p, Z}$ . If  $\psi_{\bar{S}}^{p, Y}$  is identical to  $\psi_{\bar{S}}^{p, Z}$ , then  $\psi_{\bar{S}}^{p, (Y, Z)}$  is equal to  $\psi_{\bar{S}}^{p, Y}$  and distortion will be zero. Otherwise, if  $\psi_{\bar{S}}^{p, Y}$  and  $\psi_{\bar{S}}^{p, Z}$  are different, as an active user inspects and rates products in the manner that we define, his ratings will tend to be distinguished as sampled from  $\psi_{\bar{S}}^{p, Y}$  rather than from  $\psi_{\bar{S}}^{p, Z}$ . As such, the influence of  $Z$  on predictions diminishes as  $n$  grows.

The bound depends on  $r$  through the term  $\ln(1/(1-r))$ . This term captures the dependence of KL distortion on the fraction of data produced by manipulators. As one would expect, this term vanishes when  $r$  is set to zero.

As a corollary of Theorem 1 and Proposition 1 we have the following bound on RMS distortion.

**Corollary 1.** *Fix the number of products  $N$  and let  $p$  be a linear CF algorithm. Then, for all  $M$ ,  $r \in \{0, 1/M, \dots, (M-1)/M\}$ ,  $Y \in S^{N \times (1-r)M}$ ,  $Z \in S^{N \times rM}$ , and  $\nu \in \sigma_N$ ,*

$$d_n^{\text{RMS}}(p, \nu, Y, Z) \leq \sqrt{\frac{1}{2n} \ln \frac{1}{1-r}}.$$

Figure 1 illustrates how this bound depends on  $r$  and  $n$ . The bound can offer useful guidance. For example, it ensures that if an active user has rated 22 products and no more than 10% of the training data is manipulated, then the RMS distortion induced by manipulators is less than 0.05. In a setting where users provide binary ratings and the system generates binary predictions, according to our bound on binary prediction distortion in (1) in Section 3.2, the average probability of correct predictions decreases by at most 0.05. Hence, if a binary CF system predicts ratings correctly 80% of the time in the absence of manipulation, it can maintain this probability at 75% in the presence of manipulation if it requires active users to rate 21 products before receiving recommendations.

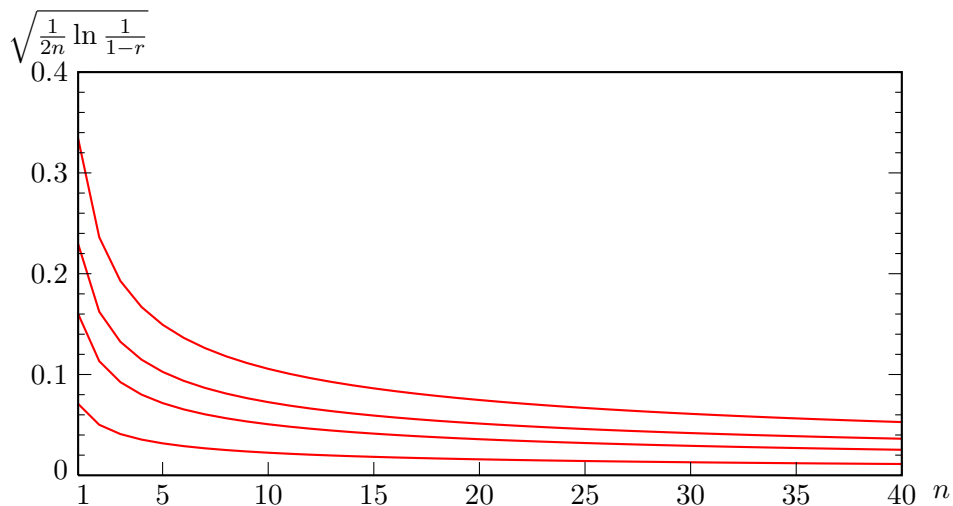


Figure 1: Bound on  $d_n^{\text{RMS}}(p, \nu, Y, Z)$  as a function of  $n$ . The four curves from bottom to top are for cases where  $r = 0.01, 0.05, 0.1,$  and  $0.2$ , respectively.

We will introduce examples of linear CF algorithms in Section 5.3.

### 4.3 Asymptotically Linear Collaborative Filtering Algorithms

We say that a probabilistic CF algorithm  $p$  is *asymptotically linear* if for all PMFs  $\psi, \phi$  over  $\bar{S}^N \times S^N$ ,  $r \in [0, 1]$ , and  $\epsilon > 0$ ,

$$\lim_{m \rightarrow \infty} \Pr \left( D \left( \left( (1-r)\hat{\psi}_{\bar{S}}^{p, U_m} + r\hat{\psi}_{\bar{S}}^{p, V_m} \right) \parallel \hat{\psi}_{\bar{S}}^{p, (U_m, V_m)} \right) \geq \epsilon \right) = 0,$$

where for each  $m$ ,  $U_m = (u^1, \dots, u^{m-l}) \in S^{N \times (m-l)}$  and  $V_m = (v^1, \dots, v^l) \in S^{N \times l}$ ,  $l \sim \text{Binomial}(m, r)$ , and  $u^1, \dots, u^{m-l} \sim \psi_S$  and  $v^1, \dots, v^l \sim \phi_S$  are i.i.d. sequences.

To understand the preceding definition, we think of training data  $(U_m, V_m)$  for each  $m$  as generated in the following way: with probability  $1 - r$ , a ratings vector is sampled from  $\psi_S$ , which we denote as  $u_i$  for an appropriate  $i$ , and with probability  $r$ , a ratings vector is sampled from  $\phi_S$ , which we denote as  $v_i$  for an appropriate  $i$ . As  $m$  grows, an asymptotically linear CF algorithm behaves like a linear CF algorithm in that the PMF  $\hat{\psi}_{\bar{S}}^{p, (U_m, V_m)}$  that it generates based on data  $(U_m, V_m)$  converges in probability to a convex combination of two PMFs: namely, the PMF  $\hat{\psi}_{\bar{S}}^{p, U_m}$  that it would generate based on the  $\psi_{\bar{S}}$ -sampled set  $U_m$  and the PMF  $\hat{\psi}_{\bar{S}}^{p, V_m}$  that it would generate based on the  $\phi_{\bar{S}}$ -sampled set  $V_m$ . By an application of the weak law of large numbers, it can be shown that all linear CF algorithms are asymptotically linear.

We can also show that asymptotically linear CF algorithms are asymptotically robust, in a sense to be made precise later. It turns out that this result applies to a broader range of practical algorithms that are asymptotically linear in a more restricted sense, which we now define. Consider a set  $\Psi$  of joint PMFs over  $\bar{S}^N \times S^N$ . We say that a probabilistic CF algorithm  $p$  is *asymptotically linear with respect to  $\Psi$*  if for all PMFs  $\psi, \phi \in \Psi$ ,  $r \in [0, 1]$ , and  $\epsilon > 0$ ,

$$\lim_{m \rightarrow \infty} \Pr \left( D \left( \left( (1-r)\hat{\psi}_{\bar{S}}^{p, U_m} + r\hat{\psi}_{\bar{S}}^{p, V_m} \right) \parallel \hat{\psi}_{\bar{S}}^{p, (U_m, V_m)} \right) \geq \epsilon \right) = 0,$$

where for each  $m$ ,  $U_m = (u^1, \dots, u^{m-l}) \in S^{N \times (m-l)}$  and  $V_m = (v^1, \dots, v^l) \in S^{N \times l}$ ,  $l \sim \text{Binomial}(m, r)$ , and  $u^1, \dots, u^{m-l} \sim \psi_S$  and  $v^1, \dots, v^l \sim \phi_S$  are i.i.d. sequences.

The following theorem and corollary characterize the robustness of asymptotically linear CF algorithms.

**Theorem 2.** *Fix the number of products  $N$  and a set  $\Psi$  of joint PMFs over  $\bar{S}^N \times S^N$ . Let  $p$  be a CF algorithm asymptotically linear with respect to  $\Psi$ . Then, for all  $\mu^*, \pi^* \in \Psi$ ,  $r \in [0, 1)$ ,  $\nu \in \sigma_N$ ,  $n \in \mathbb{Z}_+$ , and  $\epsilon > 0$ ,*

$$\lim_{m \rightarrow \infty} \Pr \left( d_n^{\text{KL}}(p, \nu, Y_m, Z_m) \geq \frac{1}{n} \ln \frac{1}{1-r} + \epsilon \right) = 0,$$

where, for each  $m$ ,  $Y_m = (y^1, \dots, y^{m-l}) \in S^{N \times (m-l)}$ ,  $Z_m = (z^1, \dots, z^l) \in S^{N \times l}$ ,  $l \sim \text{Binomial}(m, r)$ , and  $y^1, \dots, y^{m-l} \sim \mu_S^*$  and  $z^1, \dots, z^l \sim \pi_S^*$  are i.i.d. sequences.

**Corollary 2.** *Fix the number of products  $N$  and a set  $\Psi$  of joint PMFs over  $\bar{S}^N \times S^N$ . Let  $p$  be a*

CF algorithm asymptotically linear with respect to  $\Psi$ . Then, for all  $\mu^*, \pi^* \in \Psi$ ,  $r \in [0, 1)$ ,  $\nu \in \sigma_N$ ,  $n \in \mathcal{Z}_+$ , and  $\epsilon > 0$ ,

$$\lim_{m \rightarrow \infty} \Pr \left( d_n^{\text{RMS}}(p, \nu, Y_m, Z_m) \geq \sqrt{\frac{1}{2n} \ln \frac{1}{1-r}} + \epsilon \right) = 0,$$

where, for each  $m$ ,  $Y_m = (y^1, \dots, y^{m-l}) \in S^{N \times (m-l)}$ ,  $Z_m = (z^1, \dots, z^l) \in S^{N \times l}$ ,  $l \sim \text{Binomial}(m, r)$ , and  $y^1, \dots, y^{m-l} \sim \mu_S^*$  and  $z^1, \dots, z^l \sim \pi_S^*$  are i.i.d. sequences.

Theorem 2 is proved in Appendix A.2 and Corollary 2 follows from Theorem 2 and Proposition 1. These results state that for any CF algorithm  $p$  asymptotically linear with respect to  $\Psi$  and any fixed PMFs  $\mu_S^*, \pi_S^* \in \Psi$ , as honest users and manipulators sample more data from them, with high probability, the distortion bounds for linear CF algorithms in Theorem 1 and Corollary 1 will also apply to  $p$  and in particular, distortion will vanish as  $n$  grows.

The intuition behind these results is similar to that for Theorem 1. In particular, given sufficient data, the learned PMF  $\hat{\psi}_{\bar{S}}^{p, (Y_m, Z_m)}$  should closely approximate  $(1-r)\mu_S^* + r\pi_S^*$ . If  $\mu_S^*$  and  $\pi_S^*$  are similar, then  $\hat{\psi}_{\bar{S}}^{p, (Y_m, Z_m)}$  should be close to  $\mu_S^*$  and distortion should be close to zero. On the other hand, if  $\mu_S^*$  and  $\pi_S^*$  are significantly different, as an active user provides more ratings, it will be increasingly clear that they are sampled from  $\mu_S^*$  rather than  $\pi_S^*$ , and distortion will diminish.

We now study a class of asymptotically linear CF algorithms, which converge to the true PMF of user types under certain assumptions about the training data. For starters, we say a set  $\Psi$  of PMFs over  $\bar{S}^N \times S^N$  is *identifiable* if all distinct PMFs  $\psi, \phi \in \Psi$  have distinct ratings marginals  $\psi_S$  and  $\phi_S$ . Given an identifiable set  $\Psi$ , we say that a probabilistic CF algorithm  $p$  is *consistent with respect to  $\Psi$*  if for all  $\psi \in \Psi$  and  $\epsilon > 0$ ,

$$\lim_{m \rightarrow \infty} \Pr \left( D \left( \hat{\psi}_{\bar{S}}^{p, W_m} \parallel \psi_{\bar{S}} \right) \geq \epsilon \right) = 0,$$

where for each  $m$ ,  $W_m = (w_1, \dots, w_m) \in S^{N \times m}$  is generated independently and in particular,  $w_1, \dots, w_m \sim \psi_S$  is an i.i.d. sequence. This definition is meant to capture algorithms that converge to  $\psi_S$  and recover its unique corresponding joint PMF  $\psi \in \Psi$  and type marginal  $\psi_{\bar{S}}$  as the  $\psi_S$ -sampled training data grows. The following theorem states the setting in which a consistent

CF algorithm is asymptotically linear.

**Theorem 3.** *Any probabilistic CF algorithm consistent with respect to an identifiable and convex set  $\Psi$  is asymptotically linear with respect to  $\Psi$ .*

The preceding result, proved in Appendix A.2, together with the definition of consistent algorithms and Theorem 2 imply that if data are sampled i.i.d. from some PMF  $\psi$  in an identifiable and convex set  $\Psi$ , then a consistent algorithm with respect to  $\Psi$  would provide guarantees on both prediction accuracy and robustness to manipulation as training data grows. In practice, even if it is unclear whether the identifiability and convexity conditions hold, as a starting point, one might still apply a consistent CF algorithm, with the hope that it will deliver reasonable accuracy and robustness. In Section 5, we will empirically evaluate a consistent CF algorithm called the naive Bayes algorithm.

#### 4.4 Nearest Neighbor Algorithms

Nearest neighbor algorithms, widely used in commercial CF systems (Bennett, 2006; Linden et al., 2003; Ryan, 2008), generally come in two classes. The first class predicts a user’s ratings based on those provided by similar users, referred to as neighbors. The second class makes predictions on a product based on ratings that the user has provided on similar products, which can also be viewed as neighbors. In this section, we study a simple NN algorithm of the first class and the extent to which its predictions can be distorted by manipulators. We show that the bounds of the previous section do not apply to this NN algorithm, and unlike the case of linear CF algorithms, distortion does not generally diminish as the active user inspects and rates products. Though our analysis focuses on a particular NN algorithm, the resulting insights apply more broadly and in particular, to NN algorithms of the second class as well.

We study the case of binary ratings. NN algorithms identify and weight neighbors using a similarity measure. We will consider a similarity measure that increases by one for each pair of consistent ratings and decreases by one for each pair of inconsistent ratings:

$$s(x, y) = |\{1 \leq n \leq N : x_n = y_n \neq ?\}| - |\{1 \leq n \leq N : ? \neq x_n \neq y_n \neq ?\}|,$$

for any pair of ratings vectors  $x, y \in \mathcal{S}^N$ .

We consider an NN algorithm that predicts the future rating of product  $n$  for a user with ratings vector  $x$  by carrying out the following steps. First, the algorithm identifies the subset of the training data samples that offer ratings for product  $n$ . If this subset is empty, the NN algorithm optimistically predicts a rating of 1. Otherwise, from among these ratings vectors, the ones most similar to  $x$  are identified. We denote the resulting set of neighbors, which should be a singleton unless there is a tie, by  $\mathcal{N}(n, x, W)$ . Finally, an average of their ratings for product  $n$  forms the prediction:

$$\tilde{x}_{n,x,W} = \frac{\sum_{w \in \mathcal{N}(n,x,W)} w_n}{|\mathcal{N}(n,x,W)|}.$$

Our observations extend to other more complicated similarity metrics and neighbors selection methods. However, we focus on this particular case in order to keep our analysis clean.

We now consider a simple setting that facilitates analysis of RMS distortion in our NN algorithm. We are interested in how RMS distortion changes as the number of ratings  $n$  provided by an active user grows. Since  $n$  cannot exceed the number of products  $N$ , we will define an ensemble of models indexed by  $N$ . To facilitate our construction, we will only consider even  $N$ .

To keep things simple, we restrict attention to a situation where honest users agree on the ratings of all products. In particular, there is a single user type  $\bar{x}^{\text{odd}}$  which rates odd-indexed products 1 and even-indexed products 0. The user type PMF  $\mu_{\bar{S}}^*$  assigns all probability to this vector. Each honest ratings vector  $y^m$  is generated by sampling a random set of odd numbers between 1 and  $N - 1$ , then for each sample  $k$ , replacing components  $k$  and  $k + 1$  of  $\bar{x}^{\text{odd}}$  with question marks. We assume that the honest ratings  $Y$  of training data is such that each set of odd numbers between 1 and  $N - 1$  is sampled exactly once. That is, each element of  $Y$  corresponds to an element of the set  $\{(1, 0), (?, ?)\}^{N/2}$ . As such, there are  $2^{N/2}$  honest ratings vectors.

Recalling the setting that we use for assessing distortion, we now consider an active user who inspects products in the ordering  $\nu = (1, \dots, N)$ , rating each based on the prediction of the NN algorithm. It is easy to see that when there are no manipulators, the NN algorithm perfectly predicts  $\tilde{x}_{k,x^{k-1},Y} = \bar{x}_k^{\text{odd}}$ , and therefore, after the user inspects  $k$  products, his ratings history  $x^k$  has  $x_j^k = \bar{x}_j^{\text{odd}}$  for  $j \leq k$  and  $x_j^k = ?$  for  $j > k$ .

We assume that manipulators produce one half of the training data. For each honest ratings vector  $y^m$ , manipulators produce a ratings vector  $z^m$  which agrees with  $y^m$  on all products rated by  $y^m$ . However, question marks in  $y^m$  are replaced by 1 for even indices and 0 for odd indices. That is, each  $z_m$  corresponds to an element of the set  $\{(1, 0), (0, 1)\}^{N/2}$ .

Suppose  $k$  is even. Given  $x^k$ , the NN algorithm predicts what the active user's rating will be for product  $k + 1$ . To do this, it identifies neighbors  $\mathcal{N}(k + 1, x^k, (Y, Z))$ , which includes the following subsets of the training data:

- A set  $Y_1$  which consists of honest ratings vectors  $y^m$  where  $y_j^m \neq ?$  for  $j \leq k + 1$ .
- A set  $Z_1$  that, for each  $y^m \in Y_1$ , includes the corresponding manipulated vector  $z^m \in Z$ .
- A set  $Z_2$  which consists of each manipulated ratings vector  $z^m$  such that the corresponding honest ratings vector  $y^m$  has  $y_j^m \neq ?$  for  $j \leq k$  and  $y_{k+1}^m = ?$ .

Note that each of these sets is of cardinality  $2^{(N-k)/2-1}$ . Vectors in  $Y_1$  and  $Z_1$  correctly rate product  $k + 1$  as 1, whereas vectors in  $Z_2$  incorrectly rate it as 0. As a consequence, the prediction for product  $k + 1$  is  $\tilde{x}_{k+1, x^k, (Y, Z)} = 2/3$  and the resulting squared error is

$$\left(\bar{x}_{k+1}^{\text{odd}} - \tilde{x}_{k+1, x^k, (Y, Z)}\right)^2 = \frac{1}{9}.$$

The preceding argument applies for all even  $k$ . For odd  $k$ , it is easy to show that the NN algorithm correctly predicts  $\bar{x}_{k+1}^{\text{odd}} = 0$ . It follows that the RMS distortion for even  $n$  is

$$d_n^{\text{RMS}}(p, \nu, Y, Z) = \sqrt{\frac{1}{n} \sum_{k=1}^n \mathbb{E} \left[ \left( \bar{x}_k^{\text{odd}} - \tilde{x}_{k, x^{k-1}, (Y, Z)} \right)^2 \right]} = \frac{1}{3\sqrt{2}}.$$

The preceding example shows that the RMS distortion of an NN algorithm for  $r = 1/2$  does not decrease as  $n$  grows. This happens because manipulated data are strategically generated to be sufficiently similar to honest data so that no matter how many ratings an active user provides, manipulated ratings vectors will make up a fixed fraction of the neighbors and consequently induce a significant amount of distortion.

In contrast, Corollary 1 establishes that linear CF algorithms exhibit a more graceful behavior, with RMS distortion vanishing as  $n$  increases. This is not to say it is impossible to design an NN algorithm that exhibits a more desirable behavior when applied to our example. However, it is difficult to know for sure whether a given variation will behave gracefully in all relevant situations.

## 4.5 Discussion

We now provide an intuitive explanation for why linear CF algorithms should be robust to manipulation relative to NN algorithms. First note that robustness depends on how a CF algorithm learns from its mistakes. In particular, a robust algorithm should notice as it observes differences between its predictions and an active user's ratings that certain things learned from the data set are hurting rather than improving its predictions.

Recall that a linear CF algorithm  $p$  generates based on the training set  $(Y, Z)$  a PMF  $\hat{\psi}_{\bar{S}}^{p,(Y,Z)}$  that is a convex combination of  $\hat{\psi}_{\bar{S}}^{p,Y}$  and  $\hat{\psi}_{\bar{S}}^{p,Z}$ , which are PMFs that the algorithm would generate based on  $Y$  and  $Z$ , respectively. As an active user rates more products, it will be increasingly clear by probabilistic inference that his ratings  $x$  are sampled from  $\hat{\psi}_{\bar{S}}^{p,Y}$ . In effect, inaccurate predictions induced by  $Z$  will increase the weight on  $\hat{\psi}_{\bar{S}}^{p,Y}$  in the conditional PMF of ratings  $\hat{\psi}_{\bar{S}}^{p,(Y,Z)}(\cdot|x)$  conditioned on observed ratings  $x$ . And this makes future predictions more accurate.

In an NN algorithm, on the other hand, inaccurate predictions do not generally improve further predictions. In particular, manipulated ratings vectors that contribute to inaccuracies may remain in the set of neighbors while honest ratings vectors may be eliminated from it. In the example in Section 4.4, for instance, manipulated data are generated so that no matter how long an active user's ratings history is, each honest ratings vector selected as a neighbor has a manipulated counterpart that is as similar, and hence also selected as a neighbor. Consequently, as an active user provides more ratings, the numbers of honest and manipulated neighbors both decrease and stay equal. As a result, inaccurate predictions do not decrease future distortion.

## 5 Empirical Study

In this section, we present our empirical findings on the manipulation robustness of NN, linear, and asymptotically linear CF algorithms. We first introduce the data set that we worked with and then describe the methods we used to evaluate robustness.

### 5.1 Data Set

We obtained a set of movie ratings provided by users, made publicly available by Netflix’s recommendation system. Each rating is an integer between 1 and 5, which we normalized to be in  $\{0, 0.25, 0.5, 0.75, 1\}$  so that the analysis and results in our paper apply directly. We randomly sampled from the data set 5000 users, who have provided 200000 ratings of 500 movies. We then randomly chose 4000 of these users and for the purpose of our experiments, treated them as honest users and their ratings as a training set  $Y$ . We used the ratings of the other 1000 users as a test set, which we refer to as  $X$ . We then generated three separate sets of 444, 1714, and 4000 manipulated ratings vectors, respectively. Each set, which we refer to as  $Z$  for simplicity of discussion, is generated to promote 50% of the movies by using a technique reported to be effective in the literature (Burke et al., 2005; Lam and Riedl, 2004; Mehta and Nejd, 2008; Mobasher et al., 2005, 2006; O’Mahony et al., 2004; Sandvig et al., 2007; Zhang et al., 2006). Specifically, we randomly sampled 250 of the 500 movies in  $Y$  and let each manipulated ratings vector in each  $Z$  assign the highest ratings to these movies, and assign a random rating to each of the other movies, sampled from the movie’s empirical marginal PMF of ratings in  $Y$ . We then replaced a random subset of ratings in  $Z$  with question marks so that its fraction of question marks matches that in  $Y$ . Manipulated ratings vectors generated this way are meant to be similar to honest ratings vectors except on movies to promote.

### 5.2 Evaluation Methods

To test the robustness of each CF algorithm  $p$ , we treated ratings in  $X$  as ratings that an active user would provide and let  $p$  predict them. Specifically, we fixed  $n$  and for each ratings vector  $x \in X$ , identified  $n$  random products that it has assigned ratings to and randomly permuted them

to form an ordering  $\nu^x = (\nu_1^x, \dots, \nu_n^x)$ . For each  $k \leq n$ , let  $x^{k-1}$  be a ratings vector that agrees with  $x$  on products  $\nu_1^x, \dots, \nu_{k-1}^x$  and assigns question marks to the other products. An algorithm  $p$  is then used to generate a scalar prediction  $\tilde{x}_{\nu_k^x, x^{k-1}, Y}$  for the rating of product  $\nu_k^x$  based on  $x^{k-1}$  and the honest data set  $Y$ . Similarly, a prediction based on a training set  $(Y, Z)$  corrupted by manipulated ratings is denoted by  $\tilde{x}_{\nu_k^x, x^{k-1}, (Y, Z)}$ . To assess influence due to manipulation, for each  $n$ , we computed the following quantity, which we will refer to as *empirical RMS distortion*:

$$\hat{d}_n^{\text{RMS}}(p, \nu^X, X, Y, Z) = \sqrt{\frac{1}{|X|} \sum_{x \in X} \frac{1}{n} \sum_{k=1}^n \left( \tilde{x}_{\nu_k^x, x^{k-1}, Y} - \tilde{x}_{\nu_k^x, x^{k-1}, (Y, Z)} \right)^2}.$$

Here,  $\nu^X = \{(\nu_1^x, \dots, \nu_n^x) : x \in X\}$ . The empirical RMS distortion measures changes of predictions for products rated by active users. It is similar to the RMS distortion  $d_n^{\text{RMS}}(p, \nu, Y, Z)$  that we defined earlier, with one difference: whereas  $d_n^{\text{RMS}}(p, \nu, Y, Z)$  samples each  $\bar{x}_{\nu_k}$  from the PMF  $p_{\nu_k, x^{k-1}, Y}$  that the algorithm generates based on  $Y$ ,  $\hat{d}_n^{\text{RMS}}(p, \nu^X, X, Y, Z)$  uses elements of  $X$  as samples. We used empirical RMS distortion rather than RMS distortion to assess algorithms in our empirical study because computing RMS distortion would take too long, requiring a running time exponential in the number of products  $n$  rated by an active user. Further, if a CF algorithm generates a nearly correct distribution in the absence of manipulation, its empirical RMS distortion will be close to its RMS distortion.

One might also wonder whether high robustness of CF algorithms stems from high prediction accuracy or comes at the expense of it. To better understand the relationship between these two performance measures, we also computed the following RMS error for each CF algorithm, which we will refer to as *empirical RMS prediction error*:

$$\hat{\mathcal{E}}_n^{\text{RMS}}(p, \nu^X, X, Y) = \sqrt{\frac{1}{|X|} \sum_{x \in X} \frac{1}{n} \sum_{k=1}^n \left( x_{\nu_k^x} - \tilde{x}_{\nu_k^x, x^{k-1}, Y} \right)^2}.$$

This quantity computes the RMS error of predictions for ratings in  $X$  when the algorithm uses  $Y$  as training data.

For algorithms that we tested, we tuned some of their parameters by cross validation. This is a technique that selects parameter values based on the performance of the corresponding algorithm

on out-of-sample data in order to estimate their performance on future data. Specifically, we randomly sampled 20% of the users in  $Y$ . We treated their ratings as a validation set  $V$  and generated predictions based on the remaining ratings  $Y \setminus V$ . Consider a parameter  $\gamma$  that we tuned for an algorithm. For each value  $\gamma'$  in a range  $\Gamma$ , we set  $\gamma = \gamma'$ , used the corresponding algorithm to predict ratings in  $V$  based on ratings in  $Y \setminus V$ , and computed the empirical RMS prediction error  $\hat{\mathcal{E}}_n^{\text{RMS}}(p, \nu^V, V, Y \setminus V)$ . Finally, we selected a parameter value  $\gamma^*$  that results in a minimal error. Similarly, when using  $(Y, Z)$  as the training set, we sampled the validation set  $V$  from  $(Y, Z)$  and for each  $\gamma' \in \Gamma$ , computed the empirical RMS prediction error  $\hat{\mathcal{E}}_n^{\text{RMS}}(p, \nu^V, V, (Y, Z) \setminus V)$  and selected a  $\gamma^*$ . Note that we chose to optimize for prediction accuracy rather than robustness in cross validation because we wanted the algorithms to maintain reasonable accuracy and wanted to avoid tuning them to be robust for specific manipulation techniques.

Overall, for each algorithm  $p$ , we generated multiple samples of  $X$ ,  $Y$ ,  $Z$ , and  $\nu^X$ , and averaged their resultant  $\hat{d}_n^{\text{RMS}}(p, \nu^X, X, Y, Z)$  and  $\hat{\mathcal{E}}_n^{\text{RMS}}(p, \nu^X, X, Y)$  across samples to obtain reliable estimates. To summarize with our notation,  $\bar{S} = \{0, 0.25, 0.5, 0.75, 1\}$ ,  $N = 500$ ,  $(M, r) \in \{(4444, 0.1), (5714, 0.3), (8000, 0.5)\}$ , and  $1 \leq n \leq 40$ . We tested three CF algorithms: a linear CF algorithm called kernel density estimation, an asymptotically linear CF algorithm called naive Bayes, and an NN algorithm called  $k$  nearest neighbor. We now present them in detail.

### 5.3 Kernel Density Estimation Algorithms

Kernel density estimation (KDE) algorithms smooth the training data and use their resultant distribution to predict future ratings. For an in-depth treatment of KDE algorithms, see Hastie et al. (2001). In our context, we say that a probabilistic CF algorithm  $p$  is a KDE algorithm with kernels  $\{\mathcal{K}_w : w \in S^N\}$  if for any  $W \in S^{N \times M}$ ,

$$\hat{\psi}_{\bar{S}}^{p,W} = \frac{1}{M} \sum_{w \in W} \mathcal{K}_w,$$

where each  $\mathcal{K}_w$  is a PMF over  $\bar{S}^N$  parameterized by a ratings vector  $w$ . It turns out that any KDE algorithm is a linear CF algorithm and any linear CF algorithm is a KDE algorithm. We will establish this in Proposition 3 in Appendix A.3.

In our experiments, we considered a KDE algorithm with kernels  $\{\mathcal{K}_w\}$  such that for each type  $\bar{x} \in \bar{S}^N$ ,

$$\mathcal{K}_w(\bar{x}) = \prod_{n=1}^N k_{w_n}(\bar{x}_n),$$

where for each  $s \in S$ ,  $k_s$  is a PMF over  $\bar{S}$  defined as follows. For  $s \neq ?$ ,  $k_s$  is the unique PMF that satisfies  $k_s(\bar{s})/k_s(s) = \exp(-|\bar{s} - s|/\beta)$  for all  $\bar{s} \in \bar{S}$ . For  $s = ?$ ,  $k_s(\bar{s}) = 1/|\bar{S}|$  for all  $\bar{s} \in \bar{S}$ . That is,  $k_s$  assigns the highest probability to  $s$  and exponentially lower probabilities to values different from  $s$  if  $s \neq ?$ , and assigns uniform probability to all values if  $s = ?$ . It is easy to see that each  $\mathcal{K}_w$  thus defined is a PMF, and it assigns high probability to types similar to  $w$  and low probability to others. The constant  $\beta > 0$  tunes the shape of  $k_s$ , which we set to be 0.15 in our experiments.

To predict the rating of product  $\nu_n$  for a user with past ratings  $x^{n-1}$ , our KDE algorithm generates a PMF  $p_{\nu_n, x^{n-1}, W}$ , which is the conditional PMF of  $\bar{x}_{\nu_n}$  conditioned on  $x^{n-1}$  with respect to the joint PMF  $\hat{\psi}_{\bar{S}}^{p, W}$ , given by

$$p_{\nu_n, x^{n-1}, W}(\bar{s}) = \hat{\psi}_{\bar{S}}^{p, W}(\bar{x}_{\nu_n} = \bar{s} | x^{n-1}) = \frac{\sum_{w \in W} \prod_{k=1}^{n-1} k_{w_{\nu_k}}(x_{\nu_k}^{n-1}) k_{w_{\nu_n}}(\bar{s})}{\sum_{w \in W} \prod_{k=1}^{n-1} k_{w_{\nu_k}}(x_{\nu_k}^{n-1})},$$

for each  $\bar{s} \in \bar{S}$ . The corresponding scalar prediction is the expectation taken with respect to  $p_{\nu_n, x^{n-1}, W}$ :

$$\tilde{x}_{\nu_n, x^{n-1}, W} = \sum_{\bar{s} \in \bar{S}} \bar{s} p_{\nu_n, x^{n-1}, W}(\bar{s}).$$

## 5.4 Naive Bayes Algorithm

A naive Bayes (NB) algorithm assumes that the true distribution of data is a convex combination of distinct distributions in each of which features of the data are conditionally independent. It aims to learn from training data the weights of the combination and feature marginals within each distribution. For a formal analysis of the algorithm and its applications to other problem settings, see Cheeseman and Stutz (1996); Domingos and Pazzani (1997); John and Langley (1995). We now describe a particular version of the algorithm that we used and discuss the context in which it is consistent and asymptotically linear.

Our NB algorithm assumes that data are sampled from a joint PMF  $\phi$  over  $\bar{S}^N \times S^N$  such that

for each  $(\bar{w}, w) \in \bar{S}^N \times S^N$  with  $\bar{w} \rightarrow w$ ,

$$\phi(\bar{w}, w) = \left( q^{\|w\|_?} (1-q)^{N-\|w\|_?} \right) \left( \sum_{l=1}^L \eta_l \prod_{n=1}^N \theta_{l,n}^{(\bar{w}_n)} \right), \quad (2)$$

where  $q \in [0, 1]$ ,  $L \in \mathcal{Z}_+$ ,  $\eta \in T_L$ , and  $\theta_{l,n} \in T_{|\bar{S}|}$  for all  $l, n$ . Here, we write  $\bar{w} \rightarrow w$  for  $(\bar{w}, w)$  if for each  $n$ , either  $w_n = \bar{w}_n$  or  $w_n = ?$ . We let  $\|w\|_?$  denote  $|\{n : w_n = ?\}|$ . For any  $k$ , we define simplex  $T_k = \{(t_1, \dots, t_k) : t_j \geq 0, \forall 1 \leq j \leq k, \sum_{j=1}^k t_j = 1\}$ . We also let  $\theta = \{\theta_{l,n}, 1 \leq l \leq L, 1 \leq n \leq N\}$ .

To understand  $\phi$ , let us consider the following generative process of ratings vectors. Let  $\Phi = \{\phi_{\bar{S}}^1, \dots, \phi_{\bar{S}}^L\}$  be  $L$  PMFs over  $\bar{S}^N$  where each  $\phi_{\bar{S}}^l$  satisfies

$$\phi_{\bar{S}}^l(\bar{w}) = \prod_{n=1}^N \theta_{l,n}^{(\bar{w}_n)}$$

for all  $\bar{w} \in \bar{S}^N$ . That is, each  $\bar{w}_n$  is independently distributed and is equal to  $\bar{s} \in \bar{S}$  with probability  $\theta_{l,n}^{(\bar{s})}$ . A type  $\bar{w}$  is generated by first selecting a PMF from  $\Phi$  where each  $\phi_{\bar{S}}^l$  is chosen with probability  $\eta_l$  and then sampling from that PMF. A ratings vector  $w$  is then generated by randomly replacing each rating  $\bar{w}_n$  by a question mark with probability  $q$ , independent of the value  $\bar{w}_n$  and whether other ratings are replaced by question marks.

The algorithm also assumes a geometric prior for  $L$  and Dirichlet priors for  $\eta$ ,  $\theta$ , and  $q$ . Hence, the posterior probability density function (PDF) of  $(L, \eta, \theta, q)$  conditioned on training data  $W$  is given by

$$f(L, \eta, \theta, q|W) = c f(L, \eta, \theta, q) \Pr(W|L, \eta, \theta, q), \quad (3)$$

where  $c$  is a normalizing constant and prior PDF

$$f(L, \eta, \theta, q) = p_L^\tau(L) f_\eta^L(\eta) \prod_{l,n} f_\theta(\theta_{l,n}) f_q(q), \text{ with}$$

$$p_L^\tau(L) = c_L e^{-\tau L},$$

$$f_\eta^L(\eta) = c_\eta^L \prod_{l=1}^L \eta_l,$$

$$f_\theta(\theta_{l,n}) = c_\theta \prod_{\bar{s} \in \bar{S}} \theta_{l,n}^{(\bar{s})}, \forall l, n, \text{ and}$$

$$f_q(q) = c_q q(1 - q),$$

and data likelihood

$$\Pr(W|L, \eta, \theta, q) = \prod_{w \in W} \left( \left( q^{\|w\|} (1 - q)^{N - \|w\|} \right) \left( \sum_{l=1}^L \eta_l \prod_{n: w_n \neq ?} \theta_{l,n}^{(w_n)} \right) \right).$$

Here, subscripts  $L$ ,  $\eta$ ,  $\theta$ , and  $q$  of the functions  $p_L^\tau$ ,  $f_\eta^L$ ,  $f_\theta$ , and  $f_q$  denote the parameters that the distributions are over. The superscript  $\tau$  in  $p_L^\tau$  denotes dependence on parameter  $\tau$ , which controls the shape of the geometric PMF. The superscript  $L$  in  $f_\eta^L$  denotes dependence on  $L$ .  $c_L$ ,  $c_\eta^L$ ,  $c_\theta$ , and  $c_q$  are normalizing constants.

The algorithm maximizes the posterior PDF over parameters by using the expectation-maximization algorithm (Dempster et al., 1977) and obtains

$$(\hat{L}, \hat{\eta}, \hat{\theta}, \hat{q}) \in \underset{(L, \eta, \theta, q)}{\operatorname{argmax}} f(L, \eta, \theta, q|W).$$

We denote by  $\hat{\psi}_{\bar{S}}^{p,W}$  the PMF over  $\bar{S}^N$  implied by  $(\hat{L}, \hat{\eta}, \hat{\theta})$ , which the algorithm uses for predictions.

In particular, a prediction for the rating of product  $\nu_n$  for a user with past ratings  $x^{n-1}$  is given by the PMF

$$p_{\nu_n, x^{n-1}, W}(\bar{s}) = \hat{\psi}_{\bar{S}}^{p,W}(\bar{x}_{\nu_n} = \bar{s} | x^{n-1}) = \frac{\sum_{l=1}^{\hat{L}} \hat{\eta}_l \prod_{k=1}^{n-1} \hat{\theta}_{l, \nu_k}^{(x_{\nu_k}^{n-1})} \hat{\theta}_{l, \nu_n}^{(\bar{s})}}{\sum_{l=1}^{\hat{L}} \hat{\eta}_l \prod_{k=1}^{n-1} \hat{\theta}_{l, \nu_k}^{(x_{\nu_k}^{n-1})}},$$

for each  $\bar{s} \in \bar{S}$ . The corresponding scalar prediction is

$$\tilde{x}_{\nu_n, x^{n-1}, W} = \sum_{\bar{s} \in \bar{S}} \bar{s} p_{\nu_n, x^{n-1}, W}(\bar{s}).$$

In our experiments, we tuned bandwidth  $\tau$  by cross validation over the range  $\Gamma = \{1, 10, 100, 1000, 10000, 100000\}$  and settled at  $\tau = 10000$ .

We now discuss the context in which the NB algorithm is consistent and asymptotically linear.

For any  $a \in [0, 1)$ , We let  $\Psi^a$  be the set of all joint PMFs over  $\bar{S}^N \times S^N$  of the form in (2) where  $q$  is fixed to be  $a$ . We establish in Proposition 4 in Appendix A.3 that for any  $a$ ,  $\Psi^a$  is identifiable and convex, and the NB algorithm is consistent with respect to it. Then, by Theorems 2 and 3, the algorithm is asymptotically linear with respect to  $\Psi^a$  and our distortion bounds apply. Let  $\Psi_{\bar{S}}^a = \{\psi_{\bar{S}} : \psi \in \Psi^a\}$  be the set of marginals over  $\bar{S}^N$  of PMFs in  $\Psi^a$ . We note that for any  $a$ ,  $\Psi_{\bar{S}}^a$  contains all PMFs over  $\bar{S}^N$ . This implies that for any PMFs  $\mu_{\bar{S}}^*$  and  $\pi_{\bar{S}}^*$  over  $\bar{S}^N$ , if honest and manipulated data are generated by first sampling types from these PMFs and then independently replacing each rating by a question mark with the same probability  $a$ , then the NB algorithm will be asymptotically robust as the sample size grows. In our experiments, although the condition regarding replacement by question marks may not hold, we still apply the NB algorithm with the hope that it will deliver reasonable robustness.

## 5.5 $k$ Nearest Neighbor Algorithm

A class of NN algorithms called  $k$  nearest neighbor ( $k$ NN) algorithms is frequently used as a performance benchmark in prior work (Burke et al., 2005; Lam and Riedl, 2004; Zhang et al., 2006). The version that we tested works as follows.

To predict the rating of product  $\nu_n$  by a user with past ratings  $x^{n-1}$  where  $n \geq 3$ , the algorithm identifies a set of neighbors  $\mathcal{N}(\nu_n, x^{n-1}, W)$  to be  $k$  ratings vectors  $w \in W$  such that  $w_{\nu_n} \neq ?$  and score highest with  $x^{n-1}$  on the following similarity measure:

$$s(w, x^{n-1}) = \frac{\sum_{1 \leq i \leq n-1: w_{\nu_i} \neq ?} (w_{\nu_i} - \hat{w})(x_{\nu_i}^{n-1} - \hat{x}^{n-1})}{\sqrt{\sum_{1 \leq i \leq N: w_i \neq ?} (w_i - \hat{w})^2} \sqrt{\sum_{1 \leq i \leq n-1} (x_{\nu_i}^{n-1} - \hat{x}^{n-1})^2}},$$

where average ratings are given by

$$\hat{w} = \frac{\sum_{1 \leq i \leq N: w_i \neq ?} w_i}{|\{1 \leq i \leq N : w_i \neq ?\}|},$$

$$\hat{x}^{n-1} = \frac{\sum_{1 \leq i \leq n-1} x_{\nu_i}^{n-1}}{n-1}.$$

Note that  $s$  here resembles the notion of a sample correlation coefficient. Its numerator is the

covariance between non-question-mark components of  $w$  and  $x^{n-1}$ . The denominator is the product of the standard deviation of non-question-mark components of  $w$  and the same quantity for  $x^{n-1}$ . The algorithm then generates the following scalar prediction:

$$\tilde{x}_{\nu_n, x^{n-1}, W} = \min \left\{ \bar{s}_{\max}, \max \left\{ \bar{s}_{\min}, \hat{x}^{n-1} + \frac{\sum_{w \in \mathcal{N}(\nu_n, x^{n-1}, W)} s(w, x^{n-1})(w_{\nu_n} - \hat{w})}{\sum_{w \in \mathcal{N}(\nu_n, x^{n-1}, W)} |s(w, x^{n-1})|} \right\} \right\},$$

where  $\bar{s}_{\max} = \max\{\bar{s} : \bar{s} \in \bar{S}\}$  and  $\bar{s}_{\min} = \min\{\bar{s} : \bar{s} \in \bar{S}\}$ . To arrive at this quantity, for each neighbor, the difference between its rating  $w_{\nu_n}$  for product  $\nu_n$  and its average rating  $\hat{w}$  is first computed. A weighted sum of these differences is then computed, where the weights are normalized similarity measures. The user's historical ratings average  $\hat{x}^{n-1}$  is then added to the sum. The total is used as the prediction, unless it falls outside  $[\bar{s}_{\min}, \bar{s}_{\max}]$ , in which case either  $\bar{s}_{\min}$  or  $\bar{s}_{\max}$  is used, whichever is closer.

For a user with ratings history  $x^{n-1}$  where  $n \leq 2$ ,  $s(\cdot, x^{n-1})$  is not well-defined. In this case, the algorithm uses the average rating of product  $\nu_n$  in the training data to generate the prediction:

$$\tilde{x}_{\nu_n, x^{n-1}, W} = \min \left\{ \bar{s}_{\max}, \max \left\{ \bar{s}_{\min}, \frac{\sum_{\{w: w \in W, w_{\nu_n} \neq ?\}} w_{\nu_n}}{|\{w : w \in W, w_{\nu_n} \neq ?\}|} \right\} \right\}.$$

In our experiments, we tuned the number of neighbors  $k$  by cross validation over the range  $\Gamma = \{1, 2, \dots, 40\}$  and settled at  $k = 10$ .

Note that even though the  $k$ NN algorithm generates scalar predictions, it still fits our definition of CF algorithms because it is possible to come up with PMFs whose corresponding expectations equal the predictions  $\tilde{x}_{\nu_n, x^{n-1}, W}$ . We do not explicitly define such a PMF, however, because it is not necessary for computing the empirical RMS distortion in our experiments.

## 5.6 Results

Figure 2 shows the empirical RMS distortions for the three algorithms that we tested, with different fractions of manipulated data. Our results suggest that in practice, NB and KDE algorithms are significantly more robust than  $k$ NN. In particular, when a user's ratings history is short,  $k$ NN and NB both incur higher empirical RMS distortions than KDE. This difference arises because while

$k$ NN and NB ignore question marks, KDE uses them and as a result, tempers its predictions. To gain some intuition, let us consider the following problem instance where ratings are binary: the set of honest ratings  $Y$  consists of  $K$  vectors whose entries are all 1s and as many vectors whose entries are all question marks. The set of manipulated ratings  $Z$  consists of  $K$  vectors whose entries are all 0s and as many vectors whose entries are all question marks. To predict the first rating  $\bar{x}_{\nu_1}$  of an active user,  $k$ NN would yield a prediction of 1 and  $1/2$  based on  $Y$  and  $(Y, Z)$ , respectively, incurring an RMS distortion of  $1/2$ . KDE would yield a prediction close to  $3/4$  based on  $Y$  and a prediction of  $1/2$  based on  $(Y, Z)$ , incurring an RMS distortion of  $1/4$ , significantly less than that of  $k$ NN. Clearly, the presence of question marks smooths KDE’s predictions and keeps its distortion low.

In Figure 2, as more ratings are provided, distortions incurred by all three algorithms decrease. When a user’s ratings history is long, NB and KDE incur distortions significantly lower than that of  $k$ NN. Note that distortions of NB and KDE always stay below the bound in Corollary 1. The curves for  $k$ NN are flat for  $n \leq 2$  because the algorithm provides the same predictions for the first two ratings  $\bar{x}_{\nu_1}$  and  $\bar{x}_{\nu_2}$  of an active user. Note that as fraction of manipulated data  $r$  increases, distortions incurred by all three algorithms increase as well.

Figure 3 displays the empirical RMS prediction errors of the three algorithms. When  $n$  is large, their errors all decrease and in particular, NB offers the lowest error and  $k$ NN, the highest.  $k$ NN sees a spike around  $n = 3$  because the algorithm switches its prediction method there: it generates predictions by using average ratings of all users for  $n \leq 2$  and generates predictions by using average ratings of neighbors for  $n \geq 3$ .

To get a better sense of our results, we note that Netflix announced that its proprietary algorithm achieves an empirical RMS prediction error, normalized to our scale, of 0.238 on a large test set, and will award one million dollars to anyone that improves it to 0.214 (Netflix Prize, 2006). One might wonder why a decrease of 0.024 may have such a large impact on recommendation quality. We suspect that due to the large number of movies, many of them are given similar predicted ratings. As a result, a small improvement in prediction accuracy may tease apart these movies and identify the most desirable ones.

Compared to Netflix’s benchmark and target prediction errors, our results are reasonable but

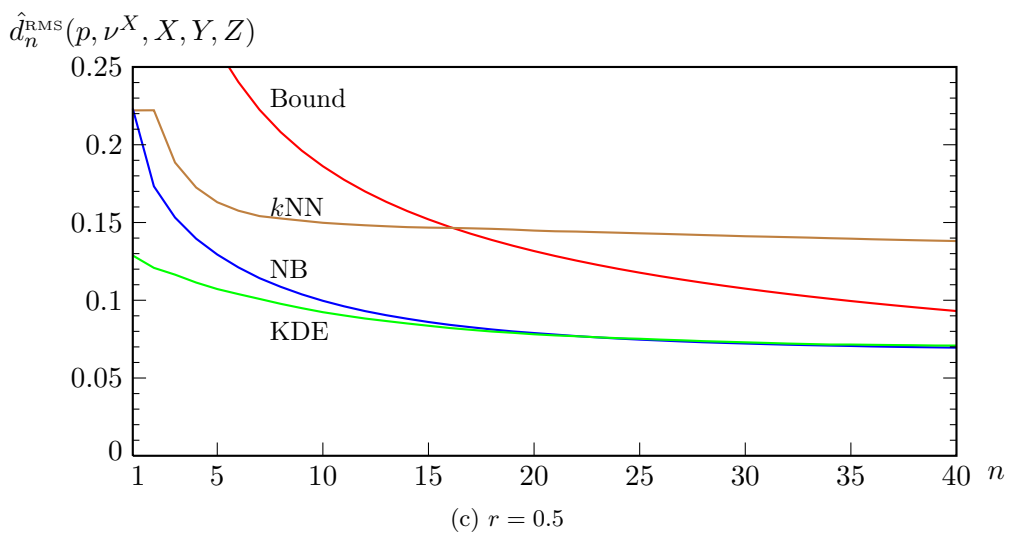
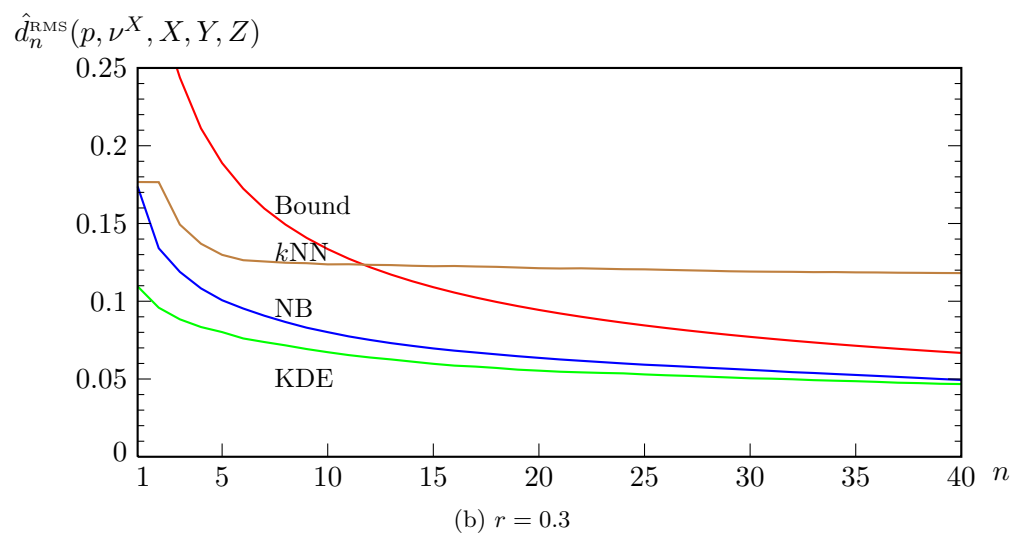
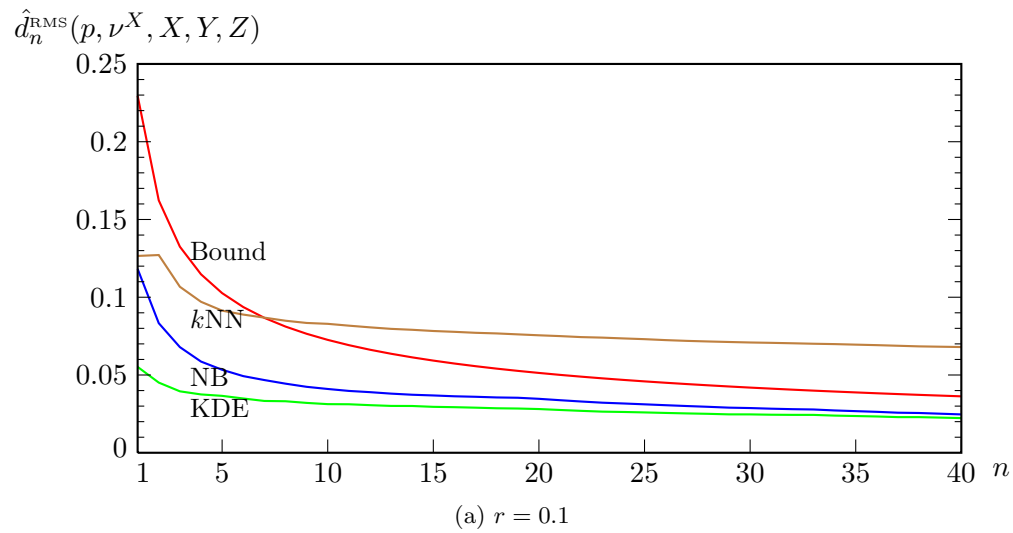


Figure 2: Empirical RMS distortion as a function of  $n$ , for different  $r$ .

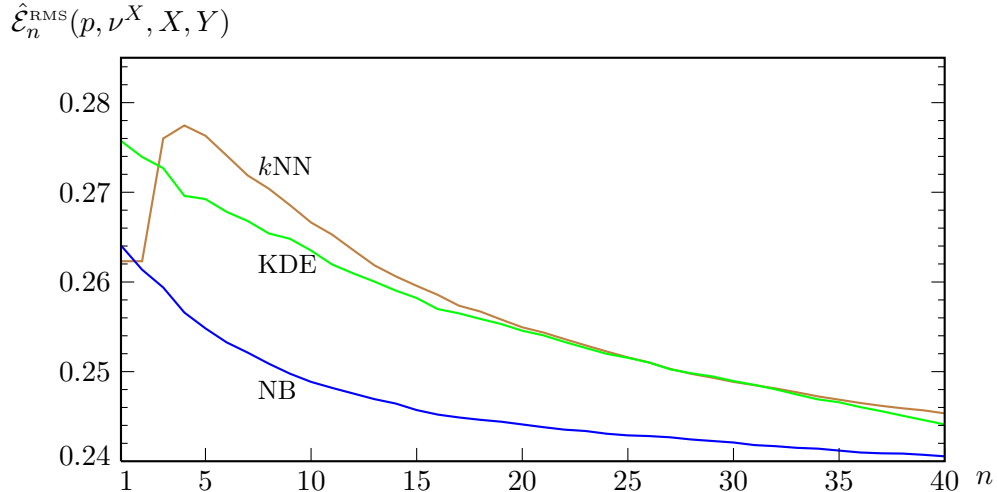


Figure 3: Empirical RMS prediction error as a function of  $n$ .

not competitive. This is because we did not focus on optimizing the prediction accuracy of the algorithms. If our objective was to achieve the highest possible accuracy while maintaining reasonable robustness, one option we could try is to fine-tune our robust algorithms to be accurate. For example, for KDE algorithms, we could work to identify more effective kernels. For NB algorithms, we could choose different priors or use methods other than expectation-maximization to find the model parameters. We could probably also design other robust linear and asymptotically linear CF algorithms that achieve higher accuracy as well. Overall, we are not suggesting that in practice, the specific algorithms that we presented should be directly implemented. Instead, one should either use them as starting points or take the insights that they yield into consideration when designing accurate and robust CF systems.

## 6 Conclusion

Our analytical and empirical work suggests that linear and asymptotically linear algorithms can be more robust to manipulation than commonly used nearest neighbor algorithms. Our results also suggest that it is possible to design algorithms that achieve accuracy alongside robustness. As such, recommendation systems of Internet commerce sites may improve their robustness to manipulation by adopting the approaches that we describe. They may also use the bounds on distortion that we

establish as a guide on how many ratings each user should provide to a recommendation system before its predictions can be trusted.

The simple setting in our work serves as a context for the initial development of our idea, and can be extended in multiple ways. One direction is to study the robustness of collaborative filtering algorithms as measured by alternative metrics. One metric could be, for instance, a user’s utility loss due to manipulation. Another extension is to design algorithms that provide non-asymptotic guarantees on both prediction accuracy and robustness.

The framework that we establish also facilitates studying the effectiveness of alternative techniques to abate influence by manipulators. For instance, given a scheme that incentivizes users to inspect and rate products, one could analyze how honest users and manipulators would behave, and then use our distortion metrics to assess the robustness of the scheme to manipulation.

It is also worth mentioning that many commercial recommendation systems build on multiple sources of information, not just collaborative filtering (Adomavicius and Tuzhilin, 2005). For example, as discussed in Balabanovic and Shoham (1997), recommendations should also be guided by features of the products being recommended. An added benefit of the approaches that we present is that they facilitate coherent fusion of multiple sources of information.

## Acknowledgments

The authors thank Christina Aperjis, Thomas Cover, Paul Cuff, Amir Dembo, Persi Diaconis, Vivek Farias, John Gill, Ramesh Johari, Yi-Hao Kao, Yi Lu, Taesup Moon, Beomsoo Park, and Assaf Zeevi for helpful suggestions. This research was supported in part by the National Science Foundation through grant IIS-0428868.

## References

- Gediminas Adomavicius and Alexander Tuzhilin. Toward the next generation of recommender systems: a survey of the state-of-the-art and possible extensions. *IEEE Transactions on Knowledge and Data Engineering*, 17(6):734–749, 2005.
- Marko Balabanovic and Yoav Shoham. Fab: Content-based, collaborative recommendation. *Communications of the ACM*, 40(3):66–72, 1997.
- James Bennett. The Cinematch system: operation, scale coverage, accuracy impact. <http://blog.recommenders06.com/wp-content/uploads/2006/09/1jimbenett.wmv>, 2006.

- Rajat Bhattacharjee and Ashish Goel. Algorithms and incentives for robust ranking. In *Proceedings of the Eighteenth Annual ACM-SIAM Symposium on Discrete Algorithms*. SIAM, 2007.
- John Breese, David Heckerman, and Carl Kadie. Empirical analysis of predictive algorithms for collaborative filtering. In *Proceedings of the Fourteenth Annual Conference on Uncertainty in Artificial Intelligence*, 1998.
- Robin Burke, Bamshad Mobasher, Roman Zabicki, and Runa Bhaumik. Limited knowledge shilling attacks in collaborative filtering systems. In *Proceedings of the Third IJCAI Workshop in Intelligent Techniques for Personalization*, 2005.
- Peter Cheeseman and John Stutz. Bayesian classification (AutoClass): Theory and results. In Usama Fayyad, Gregory Piatetsky-Shapiro, Padhraic Smyth, and Ramasamy Uthurusamy, editors, *Advances in Knowledge Discovery and Data Mining*, pages 153–180. AAAI, 1996.
- Chrysanthos Dellarocas. Strategic manipulation of internet opinion forums: Implications for consumers and firms. *Management Science*, 52(10):1577–1593, 2006.
- Arthur Dempster, Nan Laird, and Donald Rubin. Maximum likelihood from incomplete data via the EM algorithm. *Journal of the Royal Statistical Society*, 39(1):1–38, 1977.
- Pedro Domingos and Michael Pazzani. On the optimality of the simple Bayesian classifier under zero-one loss. *Machine Learning*, 29(2):103–130, 1997.
- Petros Drineas, Iordanis Kerenidis, and Prabhakar Raghavan. Competitive recommendation systems. In *Proceedings of the Thirty-fourth Annual ACM Symposium on Theory of Computing*. ACM, 2002.
- Eric Friedman, Paul Resnick, and Rahul Sami. Manipulation-resistant reputation systems. In Noam Nisan, Tim Roughgarden, Éva Tardos, and Vijay Vazirani, editors, *Algorithmic Game Theory*, chapter 27, pages 677–697. Cambridge University Press, 2007.
- Olivier Gossner and Tristan Tomala. Entropy bounds on Bayesian learning. *Journal of Mathematical Economics*, 44:24–32, 2008.
- Trevor Hastie, Robert Tibshirani, and Jerome Friedman. *The Elements of Statistical Learning*. Springer, 2001.
- Jonathan Herlocker, Joseph Konstan, Al Borchers, and John Riedl. An algorithmic framework for performing collaborative filtering. In *Proceedings of the Twenty-second Annual International ACM SIGIR Conference on Research and Development in Information Retrieval*. ACM, 1999.
- George John and Pat Langley. Estimating continuous distributions in Bayesian classifiers. In *Proceedings of the Eleventh Conference on Uncertainty in Artificial Intelligence*. Morgan Kaufmann, 1995.
- John Kleinberg and Mark Sandler. Using mixture models for collaborative filtering. *Computer and System Sciences*, 74(1):49–69, 2008.
- Shyong Lam and John Riedl. Shilling recommender systems for fun and profit. In *Proceedings of the Thirteenth International Conference on World Wide Web*. ACM, 2004.
- Greg Linden, Brent Smith, and Jeremy York. Amazon.com recommendations: Item-to-item collaborative filtering. *IEEE Internet Computing*, 7(1):76–80, 2003.
- Paolo Massa and Paolo Avesani. Trust-aware collaborative filtering for recommender systems. In *Proceedings of the Eleventh International Conference on Intelligent User Interfaces*. ACM, 2008.
- Bhaskar Mehta. Unsupervised shilling detection for collaborative filtering. In *Proceedings of the Twenty-Second Conference on Artificial Intelligence*. AAAI, 2007.

- Bhaskar Mehta and Wolfgang Nejdl. Attack resistant collaborative filtering. In *Proceedings of the Thirty-First Annual International ACM SIGIR Conference on Research and Development in Information Retrieval*. ACM, 2008.
- Nolan Miller, Paul Resnick, and Richard Zeckhauser. Eliciting informative feedback: The peer-prediction method. *Management Science*, 51(9):1359–1373, 2005.
- Bamshad Mobasher, Robin Burke, Runa Bhaumik, and Chad Williams. Effective attack models for shilling item-based collaborative filtering systems. In *Proceedings of the 2005 WebKDD Workshop*. ACM, 2005.
- Bamshad Mobasher, Robin Burke, and JJ Sandvig. Model-based collaborative filtering as a defense against profile injection attacks. In *Proceedings of the Twenty-first National Conference on Artificial Intelligence*. AAAI, 2006.
- Bamshad Mobasher, Robin Burke, Runa Bhaumik, and Chad Williams. Toward trustworthy recommender systems: An analysis of attack models and algorithm robustness. *ACM Transactions on Internet Technology*, 7(4), 2007.
- Sangkil Moon and Gary Russell. Predicting product purchase from inferred customer similarity: an autologistic model approach. *Management Science*, 54(1):71–82, 2008.
- Rajeev Motwani and Sergei Vassilvitskii. Tracing the path: new model and algorithms for collaborative filtering. In *Twenty-third International Conference on Data Engineering Workshop*. IEEE, 2007.
- Netflix Prize. <http://www.netflixprize.com>, 2006.
- John O’Donovan and Barry Smyth. Is trust robust? An analysis of trust-based recommendation. In *Proceedings of the Eleventh International Conference on Intelligent User Interfaces*. ACM, 2006.
- Stefanie Olsen. Amazon blushes over sex link gaffe. <http://news.cnet.com/2100-1023-976435.html>, 2002.
- Michael O’Mahony, Neil Hurley, Nicholas Kushmerick, and Guenole Silvestre. Collaborative recommendation: A robustness analysis. *ACM Transactions on Internet Technologies*, 4(4):344–377, 2004.
- Paul Resnick and Rahul Sami. The Influence Limiter: provably manipulation-resistant recommender systems. In *Proceedings of the ACM Recommender Systems Conference*. ACM, 2007.
- Lucas Ryan. Personal communication, 2008.
- JJ Sandvig, Bamshad Mobasher, and Robin Burke. Robustness of collaborative recommendation based on association rule mining. In *Proceedings of the 2007 ACM Conference on Recommender Systems*. ACM, 2007.
- Badrul Sarwar, George Karypis, Joseph Konstan, and John Reidl. Item-based collaborative filtering recommendation algorithms. In *Proceedings of the Tenth International Conference on World Wide Web*. ACM, 2001.
- J. Ben Schafer, Joseph Konstan, and John Riedl. E-commerce recommendation applications. *Data Mining and Knowledge Discovery*, 5(1):115–153, 2001.
- Abraham Wald. Note on the consistency of the maximum likelihood estimate. *The Annals of Mathematical Statistics*, 20(4):595–601, 1949.
- Chad Williams, Bamshad Mobasher, and Robin Burke. Defending recommender systems: Detection of profile injection attacks. *Journal of Service Oriented Computing and Applications*, 1(3), 2007.
- Sheng Zhang, Yi Ouyang, James Ford, and Fillia Makedon. Analysis of a low-dimensional linear model under recommendation attacks. In *Proceedings of the Twenty-Ninth Annual International ACM SIGIR Conference on Research and Development in Information Retrieval*. ACM, 2006.

## A Proofs

### A.1 Relationships Among Distortion Measures

Propositions 1 and 2 state relationships between KL, RMS, and binary prediction distortions. Lemmas 1 and 2 help prove them.

**Lemma 1.** *Consider two PMFs  $p$  and  $q$  with support on the same finite set  $U \subset [0, 1]$ . Let  $u$  denote a dummy variable. It holds that*

$$\left| \frac{\mathbb{E}[u]}{p} - \frac{\mathbb{E}[u]}{q} \right| \leq \frac{1}{2} \|p - q\|_1.$$

*Proof.* Let  $U = \{u_1, \dots, u_N\}$  and correspondingly, let  $p_i = p(u_i)$  and  $q_i = q(u_i)$ , for  $1 \leq i \leq N$ . Without loss of generality, let  $p_1 - q_1 \geq p_2 - q_2 \geq \dots \geq p_N - q_N$ . There exists  $n$  such that  $p_n - q_n \geq 0 \geq p_{n+1} - q_{n+1}$ . Hence,  $\sum_{i=1}^n |p_i - q_i| = \sum_{i=n+1}^N |p_i - q_i|$ . We then have

$$\begin{aligned} \left| \frac{\mathbb{E}[u]}{p} - \frac{\mathbb{E}[u]}{q} \right| &= \left| \sum_{i=1}^N u_i(p_i - q_i) \right| = \left| \sum_{i=1}^n u_i(p_i - q_i) + \sum_{i=n+1}^N u_i(p_i - q_i) \right| \\ &\leq \max \left\{ \left| \sum_{i=1}^n u_i(p_i - q_i) \right|, \left| \sum_{i=n+1}^N u_i(p_i - q_i) \right| \right\} \leq \max \left\{ \sum_{i=1}^n |p_i - q_i|, \sum_{i=n+1}^N |p_i - q_i| \right\} \\ &= \frac{1}{2} \sum_{i=1}^N |p_i - q_i|. \end{aligned}$$

□

**Proposition 1.** *Fix the number of products  $N$  and let  $p$  be a CF algorithm. Then, for all  $M$ ,  $r \in \{0, 1/M, \dots, (M-1)/M\}$ ,  $Y \in S^{N \times (1-r)M}$ ,  $Z \in S^{N \times rM}$ , and  $\nu \in \sigma_N$ ,*

$$d_n^{\text{RMS}}(p, \nu, Y, Z) \leq \sqrt{\frac{1}{2} d_n^{\text{KL}}(p, \nu, Y, Z)}.$$

*Proof.* Recall that  $\tilde{x}_{\nu_k, x^{k-1}, Y}$  and  $\tilde{x}_{\nu_k, x^{k-1}, (Y, Z)}$  denote the expected ratings of product  $\nu_k$  with

respect to PMFs  $p_{\nu_k, x^{k-1}, Y}$  and  $p_{\nu_k, x^{k-1}, (Y, Z)}$ , respectively. We have

$$\begin{aligned}
& d_n^{\text{RMS}}(p, \nu, Y, Z) \\
&= \sqrt{\frac{1}{n} \sum_{k=1}^n \mathbb{E} \left[ \left( \tilde{x}_{\nu_k, x^{k-1}, Y} - \tilde{x}_{\nu_k, x^{k-1}, (Y, Z)} \right)^2 \right]} \\
&\leq \sqrt{\frac{1}{n} \sum_{k=1}^n \mathbb{E} \left[ \left( \frac{1}{2} \| p_{\nu_k, x^{k-1}, Y} - p_{\nu_k, x^{k-1}, (Y, Z)} \|_1 \right)^2 \right]} \\
&\leq \sqrt{\frac{1}{2n} \sum_{k=1}^n \mathbb{E} \left[ D \left( p_{\nu_k, x^{k-1}, Y} \parallel p_{\nu_k, x^{k-1}, (Y, Z)} \right) \right]} \\
&= \sqrt{\frac{1}{2} d_n^{\text{KL}}(p, \nu, Y, Z)},
\end{aligned}$$

where the first inequality follows from Lemma 1 and the second inequality follows from Pinsker's inequality.  $\square$

**Lemma 2.** *Consider a Bernoulli random variable  $X$  and discrete random variables  $W_1$  and  $W_2$ . Let  $\hat{X}_1$  and  $\hat{X}_2$  be the maximum a posteriori estimates of  $X$  upon observing  $W_1$  and  $W_2$ , respectively. That is,*

$$\hat{X}_1 = \operatorname{argmax}_{x \in \{0,1\}} \Pr(X = x | W_1),$$

$$\hat{X}_2 = \operatorname{argmax}_{x \in \{0,1\}} \Pr(X = x | W_2).$$

Then,

$$\Pr(\hat{X}_1 = X) - \Pr(\hat{X}_2 = X) \leq \sqrt{\mathbb{E}[(\mathbb{E}[X|W_1] - \mathbb{E}[X|W_2])^2]}.$$

*Proof.*

$$\begin{aligned}
& \Pr(\hat{X}_1 = X) - \Pr(\hat{X}_2 = X) \\
&= \sum_{w_1} \Pr(W_1 = w_1) \max_{x \in \{0,1\}} \Pr(X = x | W_1 = w_1) - \sum_{w_2} \Pr(W_2 = w_2) \max_{x \in \{0,1\}} \Pr(X = x | W_2 = w_2) \\
&= \sum_{w_1, w_2} \Pr(W_1 = w_1, W_2 = w_2) \left( \max_x \Pr(X = x | W_1 = w_1) - \max_x \Pr(X = x | W_2 = w_2) \right) \\
&\leq \sum_{w_1, w_2} \Pr(W_1 = w_1, W_2 = w_2) |\Pr(X = 1 | W_1 = w_1) - \Pr(X = 1 | W_2 = w_2)| \\
&\leq \sqrt{\sum_{w_1, w_2} \Pr(W_1 = w_1, W_2 = w_2) (\Pr(X = 1 | W_1 = w_1) - \Pr(X = 1 | W_2 = w_2))^2} \\
&= \sqrt{\mathbb{E} \left[ (\mathbb{E}[X | W_1] - \mathbb{E}[X | W_2])^2 \right]}.
\end{aligned}$$

The first inequality follows from a simple arithmetic argument, and the second inequality follows from Jensen's inequality.  $\square$

**Proposition 2.** Fix the number of products  $N$ . Let  $p$  be a CF algorithm and  $\bar{S} = \{0, 1\}$ . Then, for all  $M, r \in \{0, 1/M, \dots, (M-1)/M\}$ ,  $Y \in S^{N \times (1-r)M}$ ,  $Z \in S^{N \times rM}$ , and  $\nu \in \sigma_N$ ,

$$d_n^{\text{B}}(p, \nu, Y, Z) \leq d_n^{\text{RMS}}(p, \nu, Y, Z).$$

*Proof.* Recall that  $\tilde{x}_{\nu_k, x^{k-1}, Y}$  and  $\tilde{x}_{\nu_k, x^{k-1}, (Y, Z)}$  denote the binary predictions on product  $\nu_k$  with respect to PMFs  $p_{\nu_k, x^{k-1}, Y}$  and  $p_{\nu_k, x^{k-1}, (Y, Z)}$ , respectively. We have

$$\begin{aligned}
& d_n^{\text{B}}(p, \nu, Y, Z) \\
&= \frac{1}{n} \sum_{k=1}^n \left( \Pr(\bar{x}_{\nu_k} = \tilde{x}_{\nu_k, x^{k-1}, Y}) - \Pr(\bar{x}_{\nu_k} = \tilde{x}_{\nu_k, x^{k-1}, (Y, Z)}) \right) \\
&\leq \sqrt{\frac{1}{n} \sum_{k=1}^n \left( \Pr(\bar{x}_{\nu_k} = \tilde{x}_{\nu_k, x^{k-1}, Y}) - \Pr(\bar{x}_{\nu_k} = \tilde{x}_{\nu_k, x^{k-1}, (Y, Z)}) \right)^2} \\
&\leq \sqrt{\frac{1}{n} \sum_{k=1}^n \mathbb{E} \left[ \left( \tilde{x}_{\nu_k, x^{k-1}, Y} - \tilde{x}_{\nu_k, x^{k-1}, (Y, Z)} \right)^2 \right]} \\
&= d_n^{\text{RMS}}(p, \nu, Y, Z).
\end{aligned}$$

The first inequality follows from Jensen's inequality and the second inequality follows from Lemma 2.  $\square$

## A.2 Results for Linear and Asymptotically Linear Collaborative Filtering Algorithms

Theorem 1 provides a distortion bound for linear CF algorithms. Theorem 2 provides a distortion bound for asymptotically linear CF algorithms. Lemmas 3 and 4 help prove it. Theorem 3 states a relationship between consistent and asymptotically CF algorithms. Lemmas 5, 6, and 7 help prove it.

**Theorem 1.** *Fix the number of products  $N$  and let  $p$  be a linear CF algorithm. Then, for all  $M$ ,  $r \in \{0, 1/M, \dots, (M-1)/M\}$ ,  $Y \in S^{N \times (1-r)M}$ ,  $Z \in S^{N \times rM}$ , and  $\nu \in \sigma_N$ ,*

$$d_n^{\text{KL}}(p, \nu, Y, Z) \leq \frac{1}{n} \ln \frac{1}{1-r}.$$

*Proof.* We denote by  $\hat{\psi}_S^{p,W,n}(\cdot|x)$  the conditional PMF of  $\bar{x}_n$  conditioned on  $x$  based on  $W$ . We have

$$\begin{aligned} & d_n^{\text{KL}}(p, \nu, Y, Z) \\ &= \frac{1}{n} \sum_{k=1}^n \mathbb{E} \left[ D \left( p_{\nu_k, x^{k-1}, Y} \parallel p_{\nu_k, x^{k-1}, (Y, Z)} \right) \right] \\ &= \frac{1}{n} \sum_{k=1}^n \mathbb{E} \left[ D \left( \hat{\psi}_S^{p, Y, \nu_k}(\cdot|x^{k-1}) \parallel \hat{\psi}_S^{p, (Y, Z), \nu_k}(\cdot|x^{k-1}) \right) \right] \\ &\leq \frac{1}{n} \sum_{k=1}^N \mathbb{E} \left[ D \left( \hat{\psi}_S^{p, Y, \nu_k}(\cdot|x^{k-1}) \parallel \hat{\psi}_S^{p, (Y, Z), \nu_k}(\cdot|x^{k-1}) \right) \right] \\ &= \frac{1}{n} D \left( \hat{\psi}_S^{p, Y} \parallel \hat{\psi}_S^{p, (Y, Z)} \right). \end{aligned}$$

The last equality follows from the chain rule of KL divergence.

For any  $\bar{x} \in \bar{S}^N$ , since  $p$  is linear, we have

$$\frac{\hat{\psi}_S^{p, Y}(\bar{x})}{\hat{\psi}_S^{p, (Y, Z)}(\bar{x})} = \frac{\hat{\psi}_S^{p, Y}(\bar{x})}{(1-r)\hat{\psi}_S^{p, Y}(\bar{x}) + r\hat{\psi}_S^{p, Z}(\bar{x})} \leq \frac{1}{1-r}.$$

Then,

$$D\left(\hat{\psi}_{\bar{S}}^{p,Y} \parallel \hat{\psi}_{\bar{S}}^{p,(Y,Z)}\right) = \sum_{\bar{x} \in \bar{S}^N} \hat{\psi}_{\bar{S}}^{p,Y}(\bar{x}) \ln \frac{\hat{\psi}_{\bar{S}}^{p,Y}(\bar{x})}{\hat{\psi}_{\bar{S}}^{p,(Y,Z)}(\bar{x})} \leq \sum_{\bar{x} \in \bar{S}^N} \hat{\psi}_{\bar{S}}^{p,Y}(\bar{x}) \ln \frac{1}{1-r} = \ln \frac{1}{1-r}.$$

Hence,

$$d_n^{\text{KL}}(p, \nu, Y, Z) \leq \frac{1}{n} D\left(\hat{\psi}_{\bar{S}}^{p,Y} \parallel \hat{\psi}_{\bar{S}}^{p,(Y,Z)}\right) \leq \frac{1}{n} \ln \frac{1}{1-r}.$$

□

**Lemma 3.** *Let  $\{\mu_m\}$  and  $\{\nu_m\}$  be two sequences of random PMFs over a fixed finite sample space  $\Omega$ . If for all  $\epsilon > 0$ ,  $\Pr(D(\mu_m \parallel \nu_m) \geq \epsilon) \rightarrow 0$ , then for all  $\epsilon > 0$ ,*

$$\Pr\left(\sum_{\omega \in \Omega} \mu_m(\omega) \left| \log \frac{\mu_m(\omega)}{\nu_m(\omega)} \right| \geq \epsilon\right) \rightarrow 0.$$

*Proof.* For  $\epsilon > 0$  and  $m$ , we denote by  $A_{m,\epsilon}$  the event that for all  $\omega \in \Omega$ , at least one of the following holds:  $|\log(\mu_m(\omega)/\nu_m(\omega))| \leq \epsilon$  and  $\max\{\mu_m(\omega), \nu_m(\omega)\} \leq \epsilon$ . We now prove that for all  $\epsilon > 0$ ,  $\Pr(A_{m,\epsilon}) \rightarrow 1$ . To see this, for any given  $\epsilon > 0$ , we let  $\delta_\epsilon$  be in  $(0, \epsilon(1 - e^{-\epsilon}))$  and denote by  $B_{m,\delta_\epsilon}$  the event that for all  $\omega \in \Omega$ ,  $|\mu_m(\omega) - \nu_m(\omega)| \leq \delta_\epsilon$ . We now show that  $B_{m,\delta_\epsilon} \subset A_{m,\epsilon}$ . If  $|\mu_m(\omega) - \nu_m(\omega)| \leq \delta_\epsilon, \forall \omega$ , then for any  $\omega'$  such that  $\max\{\mu_m(\omega'), \nu_m(\omega')\} > \epsilon$ , we have  $\min\{\mu_m(\omega'), \nu_m(\omega')\} > \epsilon - \delta_\epsilon$ . This implies

$$\begin{aligned} \left| \log \frac{\mu_m(\omega')}{\nu_m(\omega')} \right| &= \max \left\{ \log \frac{\mu_m(\omega')}{\nu_m(\omega')}, \log \frac{\nu_m(\omega')}{\mu_m(\omega')} \right\} \\ &\leq \max \left\{ \log \left( \frac{|\mu_m(\omega') - \nu_m(\omega')|}{\nu_m(\omega')} + 1 \right), \log \left( \frac{|\nu_m(\omega') - \mu_m(\omega')|}{\mu_m(\omega')} + 1 \right) \right\} < \log \left( \frac{\delta_\epsilon}{\epsilon - \delta_\epsilon} + 1 \right) \leq \epsilon, \end{aligned}$$

where the last inequality follows from our choice of  $\delta_\epsilon$ . Hence,  $B_{m,\delta_\epsilon} \subset A_{m,\epsilon}$  and for any  $\epsilon$  and a corresponding  $\delta_\epsilon$ , we have  $\Pr(A_{m,\epsilon}) \geq \Pr(B_{m,\delta_\epsilon}) \rightarrow 1$ , where convergence follows from Pinsker's inequality.

For each  $m$  and each realization of  $\mu_m$  and  $\nu_m$ , we let  $\Omega_m = \{\omega \in \Omega : \mu_m(\omega) \leq \nu_m(\omega)\}$ , let

$$\tau_m = \sum_{\omega \in \Omega_m} \mu_m(\omega) \log \frac{\mu_m(\omega)}{\nu_m(\omega)},$$

and for  $\epsilon > 0$ , denote by  $C_{m,\epsilon}$  the event that  $|\tau_m| \leq \epsilon$ . We now show that for any  $\epsilon > 0$ ,  $A_{m,\gamma_\epsilon} \subset C_{m,\epsilon}$  where  $\gamma_\epsilon \in (0, \min\{\epsilon/|\Omega|, 1/e\}]$  satisfies  $\gamma_\epsilon |\log \gamma_\epsilon| \leq \epsilon/|\Omega|$ . To see this, we first let  $\Omega_{m,\gamma_\epsilon}^1 = \{\omega \in \Omega : |\log(\mu_m(\omega)/\nu_m(\omega))| \leq \gamma_\epsilon\}$  and  $\Omega_{m,\gamma_\epsilon}^2 = \{\omega \in \Omega : \max\{\mu_m(\omega), \nu_m(\omega)\} \leq \gamma_\epsilon\} \setminus \Omega_{m,\gamma_\epsilon}^1$ . Note that  $\Omega_{m,\gamma_\epsilon}^1 \cap \Omega_{m,\gamma_\epsilon}^2 = \emptyset$ . If for all  $\omega \in \Omega$ ,  $|\log(\mu_m(\omega)/\nu_m(\omega))| \leq \gamma_\epsilon$  or  $\max\{\mu_m(\omega), \nu_m(\omega)\} \leq \gamma_\epsilon$ , then  $\Omega_{m,\gamma_\epsilon}^1 \cup \Omega_{m,\gamma_\epsilon}^2 = \Omega$ . This implies

$$\begin{aligned}
|\tau_m| &= \sum_{\omega \in \Omega_m \cap \Omega_{m,\gamma_\epsilon}^1} \mu_m(\omega) \left| \log \frac{\mu_m(\omega)}{\nu_m(\omega)} \right| + \sum_{\omega \in \Omega_m \cap \Omega_{m,\gamma_\epsilon}^2} \mu_m(\omega) \left| \log \frac{\mu_m(\omega)}{\nu_m(\omega)} \right| \\
&\leq |\Omega_{m,\gamma_\epsilon}^1| \gamma_\epsilon + \sum_{\omega \in \Omega_m \cap \Omega_{m,\gamma_\epsilon}^2} \mu_m(\omega) |\log \mu_m(\omega)| \\
&\leq |\Omega_{m,\gamma_\epsilon}^1| \gamma_\epsilon + |\Omega_{m,\gamma_\epsilon}^2| \gamma_\epsilon |\log \gamma_\epsilon| \\
&\leq |\Omega_{m,\gamma_\epsilon}^1| \frac{\epsilon}{|\Omega|} + |\Omega_{m,\gamma_\epsilon}^2| \frac{\epsilon}{|\Omega|} \\
&= \epsilon.
\end{aligned}$$

The first inequality follows from the definitions of  $\Omega_{m,\gamma_\epsilon}^1$  and  $\Omega_m$ . The second inequality follows from the definition of  $\Omega_{m,\gamma_\epsilon}^2$  and that  $\gamma_\epsilon \leq 1/e$ . The third inequality follows from other constraints on  $\gamma_\epsilon$ . Hence, for any  $\epsilon > 0$  and a corresponding  $\gamma_\epsilon$  satisfying the aforesaid constraints,  $A_{m,\gamma_\epsilon} \subset C_{m,\epsilon}$ .

Note that for all  $m$  and all realizations of  $\mu_m$  and  $\nu_m$ ,

$$\sum_{\omega \in \Omega} \mu_m(\omega) \left| \log \frac{\mu_m(\omega)}{\nu_m(\omega)} \right| = D(\mu_m \parallel \nu_m) - 2\tau_m.$$

Hence, for any  $\epsilon > 0$ , a corresponding  $\gamma_\epsilon$ , and any  $m$ ,

$$\begin{aligned}
&\Pr \left( \sum_{\omega \in \Omega} \mu_m(\omega) \left| \log \frac{\mu_m(\omega)}{\nu_m(\omega)} \right| \geq 3\epsilon \right) \\
&= \Pr(D(\mu_m \parallel \nu_m) - 2\tau_m \geq 3\epsilon, A_{m,\gamma_\epsilon}^c) + \Pr(D(\mu_m \parallel \nu_m) - 2\tau_m \geq 3\epsilon, A_{m,\gamma_\epsilon}) \\
&\leq \Pr(A_{m,\gamma_\epsilon}^c) + \Pr(D(\mu_m \parallel \nu_m) - 2\tau_m \geq 3\epsilon, C_{m,\epsilon}) \\
&\leq \Pr(A_{m,\gamma_\epsilon}^c) + \Pr(D(\mu_m \parallel \nu_m) \geq \epsilon) \rightarrow 0.
\end{aligned}$$

Here,  $A_{m,\gamma_\epsilon}^c$  denotes the complement of  $A_{m,\gamma_\epsilon}$ . The last inequality follows from the definition of  $C_{m,\epsilon}$ . Convergence follows from our original assumption and that  $\Pr(A_{m,\gamma_\epsilon}) \rightarrow 1$ .  $\square$

**Lemma 4.** Let  $\{\mu_m\}$ ,  $\{\nu_m\}$ , and  $\{\chi_m\}$  be three sequences of random PMFs over a fixed finite sample space  $\Omega$ . Suppose that for all  $\epsilon > 0$ ,  $\Pr(D(\mu_m \parallel \nu_m) \geq \epsilon) \rightarrow 0$ . Further, suppose there exists  $b > 0$  such that for all  $m$ , realizations of  $\chi_m$  and  $\mu_m$ , and  $\omega \in \Omega$ ,  $\chi_m(\omega)/\mu_m(\omega) \leq b$ . Then, for all  $\epsilon > 0$ ,  $\Pr(|D(\chi_m \parallel \mu_m) - D(\chi_m \parallel \nu_m)| \geq \epsilon) \rightarrow 0$ .

*Proof.* For any  $\epsilon > 0$  and  $m$ ,

$$\begin{aligned} & \Pr(|D(\chi_m \parallel \mu_m) - D(\chi_m \parallel \nu_m)| \geq \epsilon) \\ &= \Pr\left(\left|\sum_{\omega \in \Omega} \chi_m(\omega) \log \frac{\mu_m(\omega)}{\nu_m(\omega)}\right| \geq \epsilon\right) \\ &\leq \Pr\left(\sum_{\omega \in \Omega} \frac{\chi_m(\omega)}{\mu_m(\omega)} \left|\mu_m(\omega) \log \frac{\mu_m(\omega)}{\nu_m(\omega)}\right| \geq \epsilon\right) \\ &\leq \Pr\left(b \sum_{\omega \in \Omega} \left|\mu_m(\omega) \log \frac{\mu_m(\omega)}{\nu_m(\omega)}\right| \geq \epsilon\right) \rightarrow 0, \end{aligned}$$

where convergence follows from Lemma 3. □

**Theorem 2.** Fix the number of products  $N$  and a set  $\Psi$  of joint PMFs over  $\bar{S}^N \times S^N$ . Let  $p$  be a CF algorithm asymptotically linear with respect to  $\Psi$ . Then, for all  $\mu^*, \pi^* \in \Psi$ ,  $r \in [0, 1)$ ,  $\nu \in \sigma_N$ ,  $n \in \mathcal{Z}_+$ , and  $\epsilon > 0$ ,

$$\lim_{m \rightarrow \infty} \Pr\left(d_n^{\text{KL}}(p, \nu, Y_m, Z_m) \geq \frac{1}{n} \ln \frac{1}{1-r} + \epsilon\right) = 0,$$

where, for each  $m$ ,  $Y_m = (y^1, \dots, y^{m-l}) \in S^{N \times (m-l)}$ ,  $Z_m = (z^1, \dots, z^l) \in S^{N \times l}$ ,  $l \sim \text{Binomial}(m, r)$ , and  $y^1, \dots, y^{m-l} \sim \mu_S^*$  and  $z^1, \dots, z^l \sim \pi_S^*$  are i.i.d. sequences.

*Proof.* Let finite sample space  $\Omega = \bar{S}^N$ . For each  $m$ , let random PMFs  $\mu_m = (1-r)\hat{\psi}_{\bar{S}}^{p, Y_m} + r\hat{\psi}_{\bar{S}}^{p, Z_m}$ ,  $\nu_m = \hat{\psi}_{\bar{S}}^{p, (Y_m, Z_m)}$ , and  $\chi_m = \hat{\psi}_{\bar{S}}^{p, Y_m}$ . By the definition of asymptotically linear CF algorithms, for all  $\epsilon > 0$ ,  $\Pr(D(\mu_m \parallel \nu_m) \geq \epsilon) \rightarrow 0$ . Clearly, for all  $m$ , realizations of  $\chi_m$  and  $\mu_m$ , and  $\bar{s} \in \bar{S}^N$ ,

$\chi_m(\bar{s})/\mu_m(\bar{s}) \leq 1/(1-r)$ . Hence, for all  $\epsilon > 0$ ,

$$\begin{aligned}
& \Pr \left( d_n^{\text{KL}}(p, \nu, Y_m, Z_m) \geq \frac{1}{n} \ln \frac{1}{1-r} + \epsilon \right) \\
& \leq \Pr \left( D \left( \hat{\psi}_{\bar{S}}^{p, Y_m} \parallel \hat{\psi}_{\bar{S}}^{p, (Y_m, Z_m)} \right) \geq \ln \frac{1}{1-r} + n\epsilon \right) \\
& \leq \Pr \left( D \left( \hat{\psi}_{\bar{S}}^{p, Y_m} \parallel \hat{\psi}_{\bar{S}}^{p, (Y_m, Z_m)} \right) \geq D \left( \hat{\psi}_{\bar{S}}^{p, Y_m} \parallel (1-r)\hat{\psi}_{\bar{S}}^{p, Y_m} + r\hat{\psi}_{\bar{S}}^{p, Z_m} \right) + n\epsilon \right) \\
& \leq \Pr (|D(\chi_m \parallel \mu_m) - D(\chi_m \parallel \nu_m)| \geq n\epsilon) \rightarrow 0.
\end{aligned}$$

The second inequality holds because  $D \left( \hat{\psi}_{\bar{S}}^{p, Y_m} \parallel (1-r)\hat{\psi}_{\bar{S}}^{p, Y_m} + r\hat{\psi}_{\bar{S}}^{p, Z_m} \right) \leq \ln(1/(1-r))$  and convergence follows from Lemma 4.  $\square$

**Lemma 5.** *Let  $U \subset \mathfrak{R}^k$  be a compact set. Consider a fixed vector  $u \in U$  and a sequence of random vectors  $\{u_m\}$  for which  $\Pr(u_m \in U) \rightarrow 1$ . For any continuous function  $f : U \rightarrow \mathfrak{R}$ , if  $\Pr(\|u_m - u\|_1 \geq \epsilon) \rightarrow 0$  for all  $\epsilon > 0$ , then  $\Pr(|f(u_m) - f(u)| \geq \epsilon) \rightarrow 0$  for all  $\epsilon > 0$ .*

*Proof.* Because the continuous function  $f$  defined on compact set  $U$  is uniformly continuous, for each  $\epsilon > 0$ , there exists  $\delta > 0$  such that for any  $v, v' \in U$ ,  $\|v - v'\|_1 \leq \delta$  implies that  $|f(v) - f(v')| \leq \epsilon$ . Hence, for all  $\epsilon > 0$ ,

$$\begin{aligned}
& \Pr (|f(u_m) - f(u)| \leq \epsilon) \\
& \geq \Pr (|f(u_m) - f(u)| \leq \epsilon, u_m \in U) \\
& \geq \Pr (\|u_m - u\|_1 \leq \delta, u_m \in U) \rightarrow 1.
\end{aligned}$$

$\square$

**Lemma 6.** *Fix a finite sample space  $\Omega$ . Let  $\{\mu_m\}$  and  $\{\nu_m\}$  be two sequences of random PMFs over  $\Omega$  and let  $\mu$  be a fixed PMF over  $\Omega$ . If for all  $\epsilon > 0$ ,  $\Pr(D(\mu_m \parallel \mu) \geq \epsilon) \rightarrow 0$  and  $\Pr(D(\nu_m \parallel \mu) \geq \epsilon) \rightarrow 0$ , then for all  $\epsilon > 0$ ,  $\Pr(D(\mu_m \parallel \nu_m) \geq \epsilon) \rightarrow 0$ .*

*Proof.* We first identify the support of  $\mu$ . Without loss of generality, we let  $\mu(\omega_i) = 0, \forall 1 \leq i \leq l$  and  $\mu(\omega_i) > \delta, \forall l < i \leq |\Omega|$  for some  $l \geq 0$  and  $\delta > 0$ . In the following, we represent PMFs as vectors in  $\mathfrak{R}^{|\Omega|}$ . To this end, we define a set  $T = \{(t_1, \dots, t_{|\Omega|}) : \sum_i t_i = 1, t_i = 0, \forall i \leq l, t_i \geq \delta, \forall i > l\}$ .

Let compact set  $U = T \times T$ . We let  $u = (\mu, \mu) \in U$  and define a sequence of random vectors  $\{u_m\}$  where each  $u_m = (\mu_m, \nu_m)$ . Let continuous function  $f : U \rightarrow \mathfrak{R}$  be the KL divergence  $D(\cdot \parallel \cdot)$ . By examining the absolute continuity of  $\mu_m$  and  $\nu_m$  with respect to  $\mu$  and applying Pinsker's inequality, we have  $\Pr(u_m \in U) \rightarrow 1$  and further, for all  $\epsilon > 0$ ,  $\Pr(\|u_m - u\|_1 \geq \epsilon) \rightarrow 0$ . Hence, for all  $\epsilon > 0$ , by Lemma 5,  $\Pr(D(\mu_m \parallel \nu_m) \geq \epsilon) = \Pr(|D(\mu_m \parallel \nu_m) - D(\mu \parallel \mu)| \geq \epsilon) = \Pr(|f(u_m) - f(u)| \geq \epsilon) \rightarrow 0$ .  $\square$

**Lemma 7.** *Fix a finite sample space  $\Omega$ . Let  $\{\mu_m\}$  and  $\{\nu_m\}$  be two sequences of random PMFs over  $\Omega$  and let  $\mu$  and  $\nu$  be two fixed PMFs over  $\Omega$ . If for all  $\epsilon > 0$ ,  $\Pr(D(\mu_m \parallel \mu) \geq \epsilon) \rightarrow 0$  and  $\Pr(D(\nu_m \parallel \nu) \geq \epsilon) \rightarrow 0$ , then for all  $r \in [0, 1]$  and  $\epsilon > 0$ ,*

$$\Pr(D(((1-r)\mu_m + r\nu_m) \parallel ((1-r)\mu + r\nu)) \geq \epsilon) \rightarrow 0.$$

*Proof.* The proof here is similar to that for Lemma 6. For a fixed  $r$ , let PMF  $\chi = (1-r)\mu + r\nu$ . Without loss of generality, we let  $\chi(\omega_i) = 0, \forall 1 \leq i \leq l$  and  $\chi(\omega_i) > \delta, \forall l < i \leq |\Omega|$  for some  $l \geq 0$  and  $\delta > 0$ . In the following, we represent PMFs as vectors in  $\mathfrak{R}^{|\Omega|}$ . To this end, we define a set  $T = \{(t_1, \dots, t_{|\Omega|}) : \sum_i t_i = 1, t_i = 0, \forall i \leq l, t_i \geq \delta, \forall i > l\}$ . Let compact set  $U = T \times T$ . We let  $u = ((1-r)\mu + r\nu, (1-r)\mu + r\nu) \in U$  and define a sequence of random vectors  $\{u_m\}$  where each  $u_m = ((1-r)\mu_m + r\nu_m, (1-r)\mu + r\nu)$ . Let continuous function  $f : U \rightarrow \mathfrak{R}$  be the KL divergence  $D(\cdot \parallel \cdot)$ . By examining the absolute continuity of  $\mu_m$  with respect to  $\mu$  and that of  $\nu_m$  with respect to  $\nu$  and applying Pinsker's inequality, we have  $\Pr(u_m \in U) \rightarrow 1$ , and for all  $\epsilon > 0$ ,  $\Pr(\|u_m - u\|_1 \geq \epsilon) \rightarrow 0$ . Hence, for all  $\epsilon > 0$ , by Lemma 5,

$$\begin{aligned} & \Pr(D(((1-r)\mu_m + r\nu_m) \parallel ((1-r)\mu + r\nu)) \geq \epsilon) \\ = & \Pr(|D(((1-r)\mu_m + r\nu_m) \parallel ((1-r)\mu + r\nu)) - D(((1-r)\mu + r\nu) \parallel ((1-r)\mu + r\nu))| \geq \epsilon) \\ = & \Pr(|f(u_m) - f(u)| \geq \epsilon) \rightarrow 0. \end{aligned}$$

$\square$

**Theorem 3.** *Any probabilistic CF algorithm consistent with respect to an identifiable and convex set  $\Psi$  is asymptotically linear with respect to  $\Psi$ .*

*Proof.* We use the notation in the definition of asymptotically linear CF algorithms in Section 4.3 and Lemmas 6 and 7. Fix an algorithm  $p$  consistent with respect to  $\Psi$ . Let finite sample space  $\Omega = \bar{S}^N$ . For each  $m$ , let random PMFs  $\mu_m = \hat{\psi}_{\bar{S}}^{p,U_m}$ ,  $\nu_m = \hat{\psi}_{\bar{S}}^{p,V_m}$ , and  $\chi_m = \hat{\psi}_{\bar{S}}^{p,(U_m,V_m)}$ . Let  $\mu = \psi_{\bar{S}}$ ,  $\nu = \phi_{\bar{S}}$ . Fix  $r \in [0, 1]$ . By the consistency and convexity of  $\Psi$ , we have for all  $\epsilon > 0$ ,  $\Pr(D(\mu_m \parallel \mu) \geq \epsilon) \rightarrow 0$ ,  $\Pr(D(\nu_m \parallel \nu) \geq \epsilon) \rightarrow 0$ , and  $\Pr(D(\chi_m \parallel ((1-r)\mu + r\nu)) \geq \epsilon) \rightarrow 0$ . By Lemma 7, for all  $\epsilon > 0$ ,  $\Pr(((1-r)\mu_m + r\nu_m) \parallel ((1-r)\mu + r\nu) \geq \epsilon) \rightarrow 0$ . Then for all  $\epsilon > 0$ ,

$$\Pr\left(D\left(\left((1-r)\hat{\psi}_{\bar{S}}^{p,U_m} + r\hat{\psi}_{\bar{S}}^{p,V_m}\right) \parallel \hat{\psi}_{\bar{S}}^{p,(U_m,V_m)}\right) \geq \epsilon\right) = \Pr\left(D\left(\left((1-r)\mu_m + r\nu_m\right) \parallel \chi_m\right) \geq \epsilon\right) \rightarrow 0,$$

where convergence follows from Lemma 6.  $\square$

### A.3 Results for Kernel Density Estimation and Naive Bayes Algorithms

Propositions in this section pertain to KDE and NB algorithms.

**Proposition 3.** *Any KDE algorithm is a linear CF algorithm. Any linear CF algorithm is a KDE algorithm.*

*Proof.* Fix a KDE algorithm  $p$ . Given  $W_1 \in S^{N \times M_1}$  and  $W_2 \in S^{N \times M_2}$ , for each  $\bar{x} \in \bar{S}^N$ ,

$$\begin{aligned} & \hat{\psi}_{\bar{S}}^{p,(W_1,W_2)}(\bar{x}) \\ &= \frac{1}{M_1 + M_2} \left( \sum_{w \in W_1} \mathcal{K}_w(\bar{x}) + \sum_{w \in W_2} \mathcal{K}_w(\bar{x}) \right) \\ &= \frac{M_1}{M_1 + M_2} \left( \frac{1}{M_1} \sum_{w \in W_1} \mathcal{K}_w(\bar{x}) \right) + \frac{M_2}{M_1 + M_2} \left( \frac{1}{M_2} \sum_{w \in W_2} \mathcal{K}_w(\bar{x}) \right) \\ &= \frac{M_1}{M_1 + M_2} \hat{\psi}_{\bar{S}}^{p,W_1}(\bar{x}) + \frac{M_2}{M_1 + M_2} \hat{\psi}_{\bar{S}}^{p,W_2}(\bar{x}). \end{aligned}$$

Hence,  $p$  is linear.

To show the converse, consider a linear CF algorithm  $p'$  that generates  $\hat{\psi}_{\bar{S}}^{p',W}$  based on  $W$ . A KDE algorithm where kernel  $\mathcal{K}_w$  is set equal to  $\hat{\psi}_{\bar{S}}^{p',\{w\}}$  for each  $w \in W$  is equivalent to  $p'$ .  $\square$

**Proposition 4.** *Fix  $q \in [0, 1)$ . Let  $\Psi^q$  be the set of all joint PMFs  $\psi$  over  $\bar{S}^N \times S^N$  for which there exists  $L \in \mathcal{Z}_+$ ,  $\eta \in T_L$ , and  $\theta_{l,n} \in T_{|\bar{S}|}$  for all  $1 \leq l \leq L, 1 \leq n \leq N$  such that for each*

$(\bar{w}, w) \in \bar{S}^N \times S^N$  where  $\bar{w} \rightarrow w$ ,

$$\psi(\bar{w}, w) = \left( q^{\|w\|_?} (1-q)^{N-\|w\|_?} \right) \left( \sum_{l=1}^L \eta_l \prod_{n=1}^N \theta_{l,n}^{(\bar{w}_n)} \right).$$

Then,  $\Psi^q$  is identifiable and convex. Further, the naive Bayes algorithm is consistent with respect to  $\Psi^q$ .

*Proof.* To show that  $\Psi^q$  is identifiable, consider  $\psi, \psi' \in \Psi^q$  for which  $\psi_S = \psi'_S$ . We then have for each  $(\bar{w}, w) \in \bar{S}^N \times S^N$  where  $\bar{w} \rightarrow w$ ,

$$\begin{aligned} & \psi(\bar{w}, w) \\ &= q^{\|w\|_?} (1-q)^{N-\|w\|_?} \psi_{\bar{S}}(\bar{w}) \\ &= q^{\|w\|_?} (1-q)^{N-\|w\|_?} \psi_S(\bar{w}) (1-q)^{-N} \\ &= q^{\|w\|_?} (1-q)^{N-\|w\|_?} \psi'_S(\bar{w}) (1-q)^{-N} \\ &= q^{\|w\|_?} (1-q)^{N-\|w\|_?} \psi'_{\bar{S}}(\bar{w}) \\ &= \psi'(\bar{w}, w). \end{aligned}$$

Hence,  $\psi^q$  is identifiable.

To show that  $\Psi^q$  is convex, consider arbitrary PMFs  $\psi, \psi' \in \Psi^q$ . For each  $\lambda \in [0, 1]$ , their convex combination  $\psi^\lambda$  satisfies

$$\begin{aligned} \psi^\lambda(\bar{w}, w) &= \lambda \psi(\bar{w}, w) + (1-\lambda) \psi'(\bar{w}, w) \\ &= q^{\|w\|_?} (1-q)^{N-\|w\|_?} \left( \sum_{l=1}^L \lambda \eta_l \prod_{n=1}^N \theta_{l,n}^{(\bar{w}_n)} + \sum_{l=1}^{L'} \lambda \eta'_l \prod_{n=1}^N \theta'_{l,n}^{(\bar{w}_n)} \right) \\ &= q^{\|w\|_?} (1-q)^{N-\|w\|_?} \left( \sum_{l=1}^{L+L'} \hat{\eta}_l \prod_{n=1}^N \hat{\theta}_{l,n}^{(\bar{w}_n)} \right) \end{aligned}$$

for each  $(\bar{w}, w) \in \bar{S}^N \times S^N$  that  $\bar{w} \rightarrow w$ , where  $\hat{\eta}_l = \lambda \eta_l$  for  $l \leq L$ ,  $\hat{\eta}_l = (1-\lambda) \eta'_{l-L}$  for  $l > L$ , and for each  $n$ ,  $\hat{\theta}_{l,n} = \theta_{l,n}$  for  $l \leq L$  and  $\hat{\theta}_{l,n} = \theta_{l-L,n}$  for  $l > L$ . Hence,  $\psi^\lambda \in \Psi^q$ .  $\Psi^q$  is convex.

We now show that the naive Bayes algorithm  $p$  is consistent with respect to  $\Psi^q$  by using results in Wald (1949). We will use notation in the definition of consistent CF algorithms in Section 4.3.

We also denote by  $\psi_{\bar{S}}^*$  the true PMF over  $\bar{S}^N$  and let  $\Psi_{\bar{S}} = \{\psi_{\bar{S}} : \psi \in \Psi^q\}$  be the set of all marginals over  $\bar{S}^N$  of PMFs in  $\Psi^q$ . According to Theorem 2 in Wald (1949), if

1.  $\bar{S}^N$ ,  $\Psi_{\bar{S}}$ , and  $\psi_{\bar{S}}^*$  satisfy certain technical conditions specified in Wald (1949), and
2. There exists a constant  $b > 0$  such that for all  $m$  and all realizations of  $\{W_m\}$ ,

$$\frac{\prod_{w \in W_m} \hat{\psi}_{\bar{S}}^{p, W_m}(w)}{\prod_{w \in W_m} \psi_{\bar{S}}^*(w)} \geq b,$$

then

$$\|\hat{\psi}_{\bar{S}}^{p, W_m} - \psi_{\bar{S}}^*\|_1 \rightarrow 0, \text{ a.s.} \quad (4)$$

We verify that condition 1 holds in our problem instance and desire to find a  $b$  that satisfies condition 2. To do so, we denote by  $(L^*, \eta^*, \theta^*)$  the parameters corresponding to  $\psi_{\bar{S}}^*$  and by  $(\hat{L}^{W_m}, \hat{\eta}^{W_m}, \hat{\theta}^{W_m})$  the parameters corresponding to  $\hat{\psi}_{\bar{S}}^{p, W_m}$ . Recall that we tune the parameter  $\tau$  by cross validation over a range  $\Gamma$ . Let  $\tau^*$  be the value that we settle at. We let

$$b = \left( \min_{\tau \in \Gamma} p_L^\tau(L^*) \right) \left( f_\eta^L(\eta^*) \prod_{l,n} f_\theta(\theta_{l,n}^*) \right).$$

Because  $\hat{\psi}_{\bar{S}}^{p, W_m}$  maximizes the posterior PDF, for all  $m$  and all realizations of  $\{W_m\}$ , we have

$$\frac{\prod_{w \in W_m} \hat{\psi}_{\bar{S}}^{p, W_m}(w)}{\prod_{w \in W_m} \psi_{\bar{S}}^*(w)} \geq \frac{p_L^{\tau^*}(L^*) f_\eta^L(\eta^*) \prod_{l,n} f_\theta(\theta_{l,n}^*)}{p_L^{\tau^*}(\hat{L}^{W_m}) f_\eta^L(\hat{\eta}^{W_m}) \prod_{l,n} f_\theta(\hat{\theta}_{l,n}^{W_m})} \geq b.$$

Hence, condition 2 holds, establishing (4). By the continuity of KL divergence and properties of almost sure convergence, for all  $\epsilon > 0$ , we have

$$\Pr \left( D \left( \hat{\psi}_{\bar{S}}^{p, W_m} \parallel \psi_{\bar{S}}^* \right) \geq \epsilon \right) \rightarrow 0,$$

which implies that  $p$  is consistent with respect to  $\Psi^q$ . □

## B Table of Notation

Notation	Description
$\bar{S}$	Set of possible rating values.
$S$	Union of $\bar{S}$ and the singleton set containing the question mark.
$N$	Number of products.
$Y$	Set of ratings vectors provided by honest users.
$Z$	Set of ratings vectors provided by manipulators.
$W$	Training set consisting of all ratings vectors.
$r$	Fraction of ratings vectors generated by manipulators.
$M$	Number of ratings vectors in the training set $W$ .
$n$	Number of ratings provided by an active user.
$x^k$	Ratings vector in $S^N$ that contains $k$ ratings provided by an active user.
$\nu_k$	The index of a $k$ th inspected product by an active user.
$p_{\nu_k, x^{k-1}, W}$	PMF generated on the rating for product $\nu_k$ , for a user with history $x^{k-1}$ , based on training set $W$ .
$p$	A CF algorithm.
$\sigma_N$	Set of permutations of $\{1, \dots, N\}$ .
$d^{\text{KL}}$	KL distortion.
$d^{\text{RMS}}$	RMS distortion.
$d^{\text{B}}$	Binary prediction distortion.
$\tilde{x}_{\nu_k, x^{k-1}, W}$	Scalar prediction of rating of product $\nu_k$ for a user with history $x^{k-1}$ , based on training set $W$ .
$\tilde{x}_{\nu_k, x^{k-1}, W}$	Binary prediction of rating of product $\nu_k$ for a user with history $x^{k-1}$ , based on training set $W$ .
$\bar{w}^m$	$m$ th user type in $W$ .
$w^m$	$m$ th ratings vector in $W$ .
$\bar{w}_n$ or $\bar{x}_n$	Rating of an $n$ th product based on a user type.
$w_n$ or $x_n$	Observed rating of an $n$ th product.
$\bar{w} \rightarrow w$	For each $n$ , either $w_n = \bar{w}_n$ or $w_n = ?$ .
$\ w\ ?$	$ \{n : w_n = ?\} $ .
$\psi$	Joint PMF over $\bar{S}^N \times S^N$ .
$\psi_{\bar{S}}$	Marginal PMF over $\bar{S}^N$ .
$\psi_S$	Marginal PMF over $S^N$ .
$\mu^*$	Joint PMF over $\bar{S}^N \times S^N$ of honest users.
$\pi^*$	Joint PMF of $\bar{S}^N \times S^N$ of manipulators.
$\Psi$	Set of joint PMFs over $\bar{S}^N \times S^N$ of interest.
$\Psi_{\bar{S}}$	Set of marginal PMFs over $\bar{S}^N$ of interest.
$\hat{\psi}_{\bar{S}}^{p, W}$	PMF for $\bar{x}$ generated by algorithm $p$ based on training set $W$ .
$\hat{\psi}_{\bar{S}}^{p, W}(\cdot x)$	Conditional PMF for $\bar{x}$ conditioned on ratings vector $x \in S^N$ , generated by algorithm $p$ based on training set $W$ .
$T_k$	Simplex $\{(t_1, \dots, t_k) : t_j \geq 0, \forall 1 \leq j \leq k. \sum_{j=1}^k t_j = 1\}$ .
$\hat{d}^{\text{RMS}}$	Empirical RMS distortion.
$\hat{\mathcal{E}}^{\text{RMS}}$	Empirical RMS prediction error.
$X$	Set of ratings vectors provided by active users.
$V$	Validation set of ratings vectors.
$\nu_k^x$	The index of a $k$ th inspected product by an active user with ratings vector $x$ .

Table 1: Table of notation.