

Scalable and Sparse: Bayesian Preference Learning with Crowds

Edwin Simpson · Iryna Gurevych

Received: date

Abstract We show how to make collaborative preference learning work at scale and how it can be used to learn a target preference function from crowd-sourced data or other noisy preference labels. The collaborative model captures the reliability of each worker or data source and models their biases and error rates. It uses latent factors to share information between similar workers and a target preference function. We devise an SVI inference schema to enable the model to scale to real-world datasets. Experiments compare results using standard variational inference, laplace approximation and SVI. On real-world data we show the benefit of the personalised model over a GP preference learning approach that treats all labels as coming from the same source, as well as established alternative methods and classifier baselines. We show that the model is able to identify a number of latent features for the workers and for textual arguments.

1 Introduction

Many tasks are more suited to pairwise comparisons than classification etc. Crowds of non-expert annotators may label more accurately if presented with pairs. Implicit feedback may be taken from user actions in an application that can be represented as a preference, such as choosing an option over other options.

There are several works for learning from noisy pairwise comparisons so far (Horvitz et al. 2013 or something like that?). However, these do not provide a way to take account of item features or to model different but valid subjective viewpoints. They assume there is a single ground truth and can therefore

Ubiquitous Knowledge Processing Lab, Dept. of Computer Science, Technische Universität Darmstadt, Germany
E-mail: {simpson,gurevych}@ukp.informatik.tu-darmstadt.de

model only one task and one user’s (or a consensus of all users) preferences at once.

Work by Felt et al. 2015, Simpson et al. 2015 etc. shows that item features are particularly useful when combining crowdsourced data. A Gaussian process has not been tested for this purpose before?

GP preference learning presents a way to learn from noisy preferences but assumes constant noise and a single underlying preference function. The collaborative Gaussian process (Houlsby et al. 2012) learns multiple users’ preferences. However existing implementations do not scale and do not identify ground truth.

We show how to scale it using SVI and how to use the model to identify ground truth from subjective preferences.

In this paper, we develop methodology to solve the following questions:

1. How can we learn a rating function over large sets of items given a large number of pairwise comparisons?
2. How do we account for the different personal preferences of annotators when inferring the ground truth?

To answer these questions we make the following technical contributions:

1. We propose a method for predicting either gold-standard or personalized ratings by aggregating crowdsourced preference labels using a model of the noise and biases of individual annotators.
2. To enable this method to scale to large, real-world datasets, we develop stochastic variational inference for Bayesian matrix factorization and Gaussian process preference learning.
3. To expedite hyper-parameter tuning, we introduce a technique for gradient-based length-scale optimization of Gaussian processes.

The next section of the paper discusses related work. We then we develop our model for preference learning from crowds in Section 3, followed by our proposed inference method in Section 4 and hyper-parameter optimisation technique in Section 5. Then, in Section 6, we evaluate our approach empirically, showing first its behaviour on synthetic data, then its scalability and predictive performance on several real-world datasets.

2 Related Work

2.1 Preference Learning from Crowds

Several works have analyzed bounds on error rates or sample complexity for pairwise learning (Chen and Suh 2015; Shah et al. 2015), but do not propose methods for learning multiple rankings from crowds of users. Chen et al. (2013) account for the varying quality of pairwise labels obtained from a crowd by learning an individual model of agreement with the true pairwise labels for each worker. This approach treats the inconsistencies between annotators’ labels as noise and does not consider the items’ features. Therefore, this method

does not learn the workers’ individual preferences and cannot model how their accuracy depends on the items considered. In contrast, Fu et al. (2016) consider item features when learning to rank from pairwise labels, but do not model individual annotators at all. Uchida et al. (2017) do model the confidence of individual annotations and propose a fuzzy ranking SVM to make predictions given item features. However, their approach also assumes a single ranking over items. The benefit of jointly learning to rank and group items has also been explored (Li et al. 2018), again assuming a single ordering.

Tian and Zhu (2012) consider crowdsourcing tasks where there may be more than one correct answer. They use a nonparametric Dirichlet process model to infer a variable number of clusters of answers for each task, and also infer annotator reliability. However, they do not apply the approach to ranking using pairwise labels. Several other works learn multiple rankings from crowdsourced pairwise labels rather than a single gold-standard ranking, but do not consider the item or user features so cannot extrapolate to new users or items (Yi et al. 2013; Kim et al. 2014; Wang et al. 2016; Kim et al. 2017). Both Yi et al. (2013) and Kim et al. (2017) learn a small number of latent ranking functions that can be combined to construct personalized preferences, although neither provide a Bayesian treatment to handle data sparsity. Wang et al. (2016) consider the case where different rankings correspond to lists of items provided in response to search queries. While they model the dependence of annotator accuracy on the domain of a query, their approach was not applied to personal or subjective rankings.

A number of studies consider actively selecting pairs of items for comparison to minimize the number of pairwise labels required (Radlinski and Joachims 2007; Qian et al. 2015; Maystre and Grossglauser 2017; Cai et al. 2017). Related research treats the selection of pairwise labels as a multi-armed bandit problem (Busa-Fekete et al. 2018). In this work, we do not study the process of learning from an oracle or user that we can query. Rather, we develop a model for aggregating pairwise labels from multiple sources, which can be used as the basis of active learning methods that exploit the model uncertainty estimates provided by this Bayesian approach.

2.2 Bayesian Preference Learning

A Bayesian approach to preference learning with Gaussian processes, *GPPL*, uses item features to make predictions for unseen items and share information between similar items (Chu and Ghahramani 2005). This model assumes a single preference function over items, so cannot be used to model the individual preferences of multiple users. The approach was extended by Houlsby et al. (2012) to capture individual preferences using a latent factor model. Pairwise labels from users with common interests help to predict each other’s preference function, hence this can be seen as a *collaborative* learning method, as used in *recommender systems*. The inference techniques proposed for this model mean it scales poorly, with computational complexity $\mathcal{O}(N^3 + NP)$,

where N is the number of items and P is the number of pairwise labels, and memory complexity $\mathcal{O}(N^2 + NP + P^2)$. In this paper, we address this issue and adapt the model for aggregating crowdsourced data. An alternative to using a latent factor model is to cluster users according to preferences (Abbasnejad et al. 2013), but this is less flexible in that it does not allow for collaborative learning between users with common preferences for only subsets of items (e.g. two users may both like one genre of music, while having different preferences over other genres).

2.3 Bayesian Matrix Factorization

Preference data can be represented as an item-user matrix with N rows and M columns, where N is the number of items and M is the number of users. In this paper, we are interested in the task of predicting values in this matrix given only sparse observations of pairwise comparisons. Matrix factorization techniques are commonly used to discover latent user and item features but can fail if the data is very sparse, unless suitably regularised or given a Bayesian treatment (Salakhutdinov and Mnih 2008). Recent work on scaling Bayesian matrix factorization (BMF) to large datasets has focused on parallelizing inference (Ahn et al. 2015; Vander Aa et al. 2017; Chen et al. 2018). Instead of distributing the computation, this paper focuses on reducing the computational cost, although the method we propose is amenable to parallelization.

Several extensions of BMF use Gaussian process priors over latent factors to model correlations between items given side information or observed item features (Adams et al. 2010; Zhou et al. 2012; Houlsby et al. 2012; Bolgár and Antal 2016). However, these techniques are not directly applicable to learning from pairwise comparisons as they assume that the observations are Gaussian-distributed numerical ratings (Shi et al. 2017).

To combine Bayesian matrix factorization with a pairwise likelihood, Houlsby et al. (2012) propose a combination of expectation propagation and variational Bayesian inference. However, their proposed method does not scale sufficiently to the numbers of items, users or pairwise labels found in many important application domains. In contrast, Khan et al. (2014) develop a scalable variational EM algorithm for matrix factorization but combine this with a separate GP to model each user’s preferences. However, while the proposed method can be trained with pairwise labels, it does not capture correlations between items or users in the latent factors. Furthermore, their scalable inference method sub-samples training data rather than learning from the complete training set.

2.4 Stochastic Variational Inference

Models that combine Gaussian processes with non-Gaussian likelihoods require approximate inference methods that often scale poorly with the amount

of training data available. This problem can be tackled using *Stochastic variational inference (SVI)* (Hoffman et al. 2013). SVI has been successfully applied to Gaussian processes (Hensman et al. 2013), including Gaussian process classifiers (Hensman et al. 2015), and Gaussian process preference learning (GPPL) (Simpson and Gurevych 2018). This paper adapts SVI to Bayesian matrix factorization for the first time as part of a solution for collaborative preference learning. We also provide the first full derivation of SVI for GPPL and introduce a technique for efficiently tuning the length-scale hyperparameters of the Gaussian processes.

3 Bayesian Preference Learning for Crowds

3.1 Modeling Pairwise Preferences

A pairwise comparison, $y(a \succ b)$, between items a and b has a binary label that is one if a is preferred to b , or zero if b is preferred to a (also written $a \prec b$). We assume that the likelihood of pairwise label $y(a, b)$ depends on the underlying value of the items to the user, represented through a latent function of the items' features, $f(\mathbf{x}_a)$, where \mathbf{x}_b is a vector representation of the features of item a . The relationship between the value function, f , and the pairwise labels can be modeled by any of several different likelihood functions, including the Bradley-Terry model (Bradley and Terry 1952; Plackett 1975; Luce 1959) and the Thurstone-Mosteller model (Thurstone 1927; Mosteller 2006). The Bradley-Terry model takes the following form:

$$p(y(a \succ b)|f) = \frac{1}{1 + \exp(f(\mathbf{x}_a) - f(\mathbf{x}_b))} \quad (1)$$

This is a logistic likelihood, which allows pairwise labels that do not reflect the true relative values of the items, due to labeling errors, variability in the user's judgements, or if the preferences are derived from noisy implicit data such as clicks streams. The error rate is determined by the relative difference in f values of the items.

A different view is to treat the errors as the result of random noise in the value function:

$$p(y(a \succ b)|f, \delta_a, \delta_b) = \begin{cases} 1 & \text{if } f(\mathbf{x}_a) + \delta_a \geq f(\mathbf{x}_b) + \delta_b \\ 0 & \text{otherwise,} \end{cases} \quad (2)$$

where $\delta \sim \mathcal{N}(0, 0.5)$ is Gaussian-distributed noise. Integrating out the unknown values of δ_a and δ_b , we get a probit likelihood:

$$\begin{aligned} p(y(a \succ b)|f) &= \int \int p(y(a \succ b)|f, \delta_a, \delta_b) \mathcal{N}(\delta_a; 0, 0.5) \mathcal{N}(\delta_b; 0, 0.5) d\delta_a d\delta_b \\ &= \Phi(z), \end{aligned} \quad (3)$$

where $z = f(\mathbf{x}_a) - f(\mathbf{x}_b)$, and Φ is the cumulative distribution function of the standard normal distribution. This is a Thurstone-Mosteller model, sometimes referred to as *Thurstone case V*, and was used for Gaussian process preference learning (GPPL) by Chu and Ghahramani (2005), with the difference that they learned the variance of the random noise, δ rather than assuming it is 0.5. However, this is unnecessary in practice, since we scale instead the value function, f , to reduce or increase the certainty in the pairwise labels. Both the logistic and probit approaches can be used here, but we proceed with the probit likelihood (as in (Herbrich et al. 2007; Chu and Ghahramani 2005)) because it allows us to handle uncertainty in f in a simple manner by modifying z :

$$\hat{z} = \frac{\mu_a - \mu_b}{\sqrt{1 + \sigma_a + \sigma_b - \sigma_{a,b}}} \quad (4)$$

where μ_a and μ_b are the expected values of $f(\mathbf{x}_a)$ and $f(\mathbf{x}_b)$ respectively, σ_a and σ_b are the corresponding variances, and $\sigma_{a,b}$ is the covariance between $f(\mathbf{x}_a)$ and $f(\mathbf{x}_b)$.

3.2 Single User Preference Learning

First consider modeling the preferences of a single user. In this case, we assume that the value function, f , is a function of item features and has a Gaussian process prior: $f \sim \mathcal{GP}(0, k_\theta/s)$, where k_θ is a kernel function with hyper-parameters θ , and $1/s$ is the scale of the function drawn from a gamma prior, $s \sim \mathcal{G}(\alpha_0, \beta_0)$, with shape α_0 and scale β_0 . The value of s determines the variance of f and therefore its magnitude, which affects the level of certainty in the pairwise label likelihood, Equation 3. The kernel function takes item features as inputs and determines the covariance between values of f for different items. Typically, we choose a kernel function that produces higher covariance between items with similar feature values, such as the *squared exponential* or *Matérn* functions. The choice of kernel function is a model selection problem as it controls the shape and smoothness of the function across the feature space. However, the Matérn and squared exponential make minimal assumptions and so are effective in a wide range of tasks (see Rasmussen and Williams (2006) for more).

We observe a set of P pairwise preference labels for a single user, $\mathbf{y} = \{y_1, \dots, y_P\}$, where the p th label, $y_p = y(a_p \succ b_p)$. The joint distribution over all variables is as follows:

$$\begin{aligned} p(\mathbf{y}, \mathbf{f}, s | k_\theta, \alpha_0, \beta_0) &= \prod_{p=1}^P p(y_p | \mathbf{f}) \mathcal{N}(\mathbf{f}; \mathbf{0}, \mathbf{K}_\theta/s) \mathcal{G}(s; \alpha_0, \beta_0) \\ &= \prod_{p=1}^P \Phi(z_p) \mathcal{N}(\mathbf{f}; \mathbf{0}, \mathbf{K}_\theta/s) \mathcal{G}(s; \alpha_0, \beta_0), \end{aligned} \quad (5)$$

where $\mathbf{f} = \{f(\mathbf{x}_1), \dots, f(\mathbf{x}_N)\}$ are the latent values for the N items referred to by the pairwise labels \mathbf{y} , and θ , α_0 and β_0 are hyper-parameters.

3.3 Latent components: Bayesian Matrix Factorization

We wish to exploit similarities between the value functions of different users or label sources to improve our preference model, particularly when faced with sparse data. In a scenario with multiple users or label sources, we can represent preference values in a matrix, \mathbf{F} , where rows correspond to items, columns to users, and entries are preference values. If we factorize this matrix, we obtain two lower-dimensional matrices, one for users, $\mathbf{W} \in \mathcal{R}^{C \times U}$, and one for the items, $\mathbf{V} \in \mathcal{R}^{N \times C}$, where C is the number of latent components, U is the number of users, and N is the number of items: $\mathbf{F} = \mathbf{V}^T \mathbf{W}$. Each row of \mathbf{V} matrices is a vector representation of an item, while each row of \mathbf{W} is a vector representation of a user, both containing the values of latent features. Users with similar values for a certain feature will have similar preferences for the subset of items with corresponding feature values. The features could represent, for example, in the case of book recommendation, interests in a particular genre of book. Using vector representations for users and items reflects that users may have overlapping sets of interests, and that items may have multiple features that make them attractive.

Besides latent features, we may also observe a number of item features, \mathbf{x} , and user features, \mathbf{u} . In the single user model, we assumed a single latent value function, f , of the observed item features. For the multi-user case, we assume that there are C latent functions, v_c over item features and C latent functions, w_c , over user features, and thereby model the relationship between each observed feature and each of the latent features. The matrices \mathbf{V} and \mathbf{W} are evaluations of these functions at the points corresponding to the observed users and items. Therefore, the latent preference function, f , for a user with features \mathbf{u} is a weighted sum over latent functions:

$$f(\mathbf{x}_a, \mathbf{u}_j) = \sum_{c=1}^C w_c(\mathbf{u}_j) v_c(\mathbf{x}_a) \quad (6)$$

To provide a Bayesian treatment to matrix factorization, we place Gaussian process priors over the latent functions:

$$v_c \sim \mathcal{GP}(\mathbf{0}, k_\theta / s_c) \quad w_c \sim \mathcal{GP}(\mathbf{0}, k_\theta). \quad (7)$$

It is not necessary to learn a separate scale for w_c , since v_c and w_c are multiplied with each other, making a single s_c equivalent to the product of two separate scales. The choice of C can be treated as a hyperparameter, or modeled using a non-parametric prior, such as the Indian Buffet Process, which assumes an infinite number of latent components (Ding et al. 2010). For simplicity, we assume fixed values of C in this paper, and allow the scale parameter $s_c \approx 0$ to effectively remove any dimensions that are not required to model the data. This section described a Bayesian matrix factorization model, which we will subsequently extend to a preference learning model for crowds of users and label sources.

3.4 Crowd Preference Learning

We combine the matrix factorization method with the preference likelihood of Equation 3 to obtain a joint preference model for multiple users or label sources. In addition to the latent components, we introduce a common value function over item features, $t \sim \mathcal{GP}(\mathbf{0}, k_\theta/\sigma_t)$, that is shared across all users. Its values $\mathbf{t} = \{t(\mathbf{x}_1), \dots, t(\mathbf{x}_N)\}$ represent a consensus between users, if present, while allowing individual users' preferences to deviate from this value through $\mathbf{V}^T \mathbf{W}$. Hence, \mathbf{t} can model the underlying ground truth or consensus in crowd-sourcing scenarios, or when using multiple label sources to learn preferences for one individual. The joint distribution of this crowd model is:

$$p(\mathbf{y}, \mathbf{V}, \mathbf{W}, \mathbf{t}, s_1, \dots, s_C, \sigma_t | k_\theta, \alpha_0, \beta_0) = \prod_{p=1}^P \Phi(z_p) \mathcal{N}(\mathbf{t}; \mathbf{0}, \mathbf{K}_{t,\theta}/\sigma_t) \mathcal{G}(\sigma_t; \alpha_0, \beta_0) \\ \prod_{c=1}^C \{\mathcal{N}(\mathbf{v}_c; \mathbf{0}, \mathbf{K}_{v,\theta}/s_c) \mathcal{N}(\mathbf{w}_c; \mathbf{0}, \mathbf{K}_{w,\theta}) \mathcal{G}(s_c; \alpha_0, \beta_0)\}, \quad (8)$$

$$\text{where } z_p = \mathbf{v}_{\cdot, a_p}^T \mathbf{w}_{\cdot, u_p} + t_{a_p} - \mathbf{v}_{\cdot, b_p}^T \mathbf{w}_{\cdot, u_p} - t_{b_p}, \quad (9)$$

and σ_t is the inverse scale of t . The index p , which identifies one observation, now refers to a tuple, $\{u_p, a_p, b_p\}$ that identifies the user and a pair of items.

4 Scalable Inference

Given a set of pairwise labels, \mathbf{y} , the goal is to infer the posterior distribution over the preference values \mathbf{f} , in the single user case, or $\mathbf{F} = \mathbf{V}^T \mathbf{W}$ in the multi-user case. Previous approaches include a Laplace approximation for the single user case (Chu and Ghahramani 2005) and a combination of expectation propagation (EP) with variational Bayes (VB) for a multi-user model (Houlsby et al. 2012). The Laplace approximation is a maximum a-posteriori (MAP) solution that takes the most probable values of parameters rather than integrating over their distributions and has been shown to perform poorly for tasks such as classification (Nickisch and Rasmussen 2008). EP approximates the true posterior with a simpler, factorized distribution that can be learned using an iterative algorithm. The true posterior is multi-modal, since the latent factors can be re-ordered arbitrarily without affecting \mathbf{F} : this is the non-identifiability problem. A standard EP approximation would average these modes before predicting \mathbf{F} , producing uninformative predictions over \mathbf{F} . Houlsby et al. (2012) resolve this by incorporating a VB step, which approximates a single mode. A drawback of EP is that unlike VB, convergence is not guaranteed (Minka 2001).

The cost of inference can be reduced using a *sparse* approximation based on a set of *inducing points*, which act as substitutes for the set of points in

the training dataset. By choosing a fixed number of inducing points, $M \ll N$, the computational cost is fixed at $\mathcal{O}(M^3)$. These points must be selected so as to give a good approximation, using either heuristics or optimizing their positions to maximize the approximate marginal likelihood. Houlsby et al. (2012) use a FITC approximation (Snelson and Ghahramani 2006) with their EP method to limit the costs of inference. However, in practice, FITC is unsuitable for datasets with more than a few thousands points as it is not amenable to distributed computation, does it address other expensive operations with computational complexity $\mathcal{O}(NP)$ and memory complexity $\mathcal{O}(P^2 + NP + N^2)$, which may become limiting when the number of data points is large (Hensman et al. 2015). We turn to stochastic variational inference (SVI) (Hoffman et al. 2013) to derive a more scalable approach for Gaussian process preference learning, including a multi-user model founded on Bayesian matrix factorization. First, we define an approximate posterior that can be estimated using SVI, then provide the update equations for an iterative algorithm to optimize this approximation. We begin with the model for a single user, then extend this to the multi-user case using matrix factorization.

4.1 An Approximate Preference Likelihood

Due to the non-Gaussian likelihood, Equation 3, the posterior distribution over \mathbf{f} contains intractable integrals:

$$p(\mathbf{f}|\mathbf{y}, k_\theta, \alpha_0, \beta_0) = \frac{\int \prod_{p=1}^P \Phi(z_p) \mathcal{N}(\mathbf{f}; \mathbf{0}, \mathbf{K}_\theta/s) \mathcal{G}(s; \alpha_0, \beta_0) ds}{\int \int \prod_{p=1}^P \Phi(z_p) \mathcal{N}(\mathbf{f}'; \mathbf{0}, \mathbf{K}_\theta/s) \mathcal{G}(s; \alpha_0, \beta_0) ds d\mathbf{f}'}. \quad (10)$$

To simplify the integral in the denominator, we approximate the preference likelihood with a Gaussian:

$$\prod_{p=1}^P \Phi(z_p) \approx \mathcal{N}(\mathbf{y}; \Phi(\mathbf{z}), \mathbf{Q}), \quad (11)$$

where $\mathbf{z} = \{z_1, \dots, z_P\}$ and \mathbf{Q} is a diagonal noise covariance matrix. We estimate the diagonal entries of \mathbf{Q} by moment matching the approximate likelihood with a beta-binomial with variance given by:

$$Q_{p,p} = \mathbb{E}_{\mathbf{f}}[\Phi(z_p)(1 - \Phi(z_p))] = \frac{(y_p + \gamma_0)(1 - y_p + \lambda_0)}{(2 + \gamma_0 + \lambda_0)}, \quad (12)$$

where γ_0 and λ_0 are parameters of a Bernoulli distribution that has the same variance as the prior $p(\Phi(z_p)|\mathbf{K}_\theta, \alpha_0, \beta_0)$ using numerical integration. Setting \mathbf{Q} in this way matches the moments of the true likelihood, $\Phi(z_p)$, to those of the Gaussian approximation.

Unfortunately, the nonlinear term, $\Phi(\mathbf{z})$ means that the posterior is still intractable, so we linearize $\Phi(\mathbf{z})$ by taking its first-order Taylor series expansion

about the expected value of \mathbf{f} :

$$\Phi(\mathbf{z}) \approx \tilde{\Phi}(\mathbf{z}) = \mathbf{G}(\mathbf{f} - \mathbb{E}[\mathbf{f}]) + \Phi(\mathbb{E}[\mathbf{z}]), \quad (13)$$

$$G_{p,i} = \Phi(\mathbb{E}[z_p])(1 - \Phi(\mathbb{E}[z_p]))(2y_p - 1)([i = a_p] - [i = b_p]) \quad (14)$$

where \mathbf{G} is a matrix containing elements $G_{p,i}$, which are the partial derivatives of the pairwise likelihood with respect to each of the latent function values, \mathbf{f} . This creates a dependency between the posterior mean of \mathbf{f} and the linearization terms in the likelihood, which can be estimated iteratively using variational inference (Steinberg and Bonilla 2014), as we will describe below. The linearization makes the approximate likelihood conjugate to $\mathcal{N}(\mathbf{f}; \mathbf{0}, \mathbf{K}_\theta/s)$, so that the approximate posterior over \mathbf{f} is also Gaussian.

Given our likelihood approximation, we can now use variational inference to estimate the marginal over \mathbf{f} , by optimizing an approximate posterior over all latent variables:

$$\begin{aligned} p(\mathbf{f}, s | \mathbf{y}, \mathbf{x}, k_\theta, \alpha_0, \beta_0) &\approx q(s)q(\mathbf{f}), \\ \text{where } \log q(s) &= \log \mathcal{N}(\mathbb{E}[\mathbf{f}]; \mathbf{0}, \mathbf{K}_\theta/s) + \log \mathcal{G}(s; \alpha_0, \beta_0) + \text{const}, \\ \log q(\mathbf{f}) &= \log \mathcal{N}(\mathbf{y}; \tilde{\Phi}(\mathbf{z}), \mathbf{Q}) + \log \mathcal{N}(\mathbf{f}; \mathbf{0}, \mathbf{K}_\theta/\mathbb{E}[s]) + \text{const}. \end{aligned} \quad (15)$$

The Gaussian likelihood approximation and linearization also appear in methods based on expectation propagation (Rasmussen and Williams 2006) and the extended Kalman filter (Reece et al. 2011; Steinberg and Bonilla 2014). However, neither these methods nor our approximation in Equation 15 make inference sufficiently scalable, as they all require inversion of an $N \times N$ matrix and further computations involving $N \times P$ and $P \times P$ matrices. We therefore modify Equation 15 to enable stochastic variational inference (SVI).

4.2 Sparse Approximation for the Single-User Model

We introduce a sparse approximation to the Gaussian process that allows us to limit the size of the covariance matrix that needs to be inverted, and permit stochastic inference methods that consider only a subset of the P observations at each iteration (Hensman et al. 2013, 2015). To do this, we introduce a set of M *inducing points*, with inputs \mathbf{x}_m , function values \mathbf{f}_m , and covariance \mathbf{K}_{mm} . The covariance between the observations and the inducing points is \mathbf{K}_{nm} . We then modify the variational approximation in Equation 15 to introduce the inducing points (for clarity, we omit θ from this point on):

$$p(\mathbf{f}, \mathbf{f}_m, s | \mathbf{y}, \mathbf{x}, \mathbf{x}_m, k_\theta, \alpha_0, \beta_0) \approx q(\mathbf{f}, \mathbf{f}_m, s) = q(s)q(\mathbf{f})q(\mathbf{f}_m), \quad (16)$$

$$\begin{aligned} \log q(\mathbf{f}_m) &= \log \mathcal{N}(\mathbf{y}; \tilde{\Phi}(\mathbf{z}), \mathbf{Q}) + \log \mathcal{N}(\mathbf{f}_m; \mathbf{0}, \mathbf{K}_{mm}/\mathbb{E}[s]) + \text{const}, \\ &= \log \mathcal{N}(\mathbf{f}_m; \hat{\mathbf{f}}_m, \mathbf{S}), \end{aligned} \quad (17)$$

$$\mathbf{S}^{-1} = \mathbf{K}_{mm}^{-1}/\mathbb{E}[s] + \mathbf{A}^T \mathbf{G}^T \mathbf{Q}^{-1} \mathbf{G} \mathbf{A}, \quad (18)$$

$$\hat{\mathbf{f}}_m = \mathbf{S} \mathbf{A}^T \mathbf{G}^T \mathbf{Q}^{-1} (\mathbf{y} - \Phi(\mathbb{E}[\mathbf{z}]) + \mathbf{G} \mathbb{E}[\mathbf{f}]), \quad (19)$$

where $\mathbf{A} = \mathbf{K}_{nm}\mathbf{K}_{mm}^{-1}$. The factor $q(s)$ remains unchanged from Equation 15, while $q(\mathbf{f})$ is now assumed to be independent of the observations:

$$\log q(\mathbf{f}) = \log \mathcal{N}(\mathbf{f}; \mathbf{A}\hat{\mathbf{f}}_m, \mathbf{K} + \mathbf{A}(\mathbf{S} - \mathbf{K}_{mm}/\mathbb{E}[s])\mathbf{A}^T). \quad (20)$$

The use of inducing points therefore avoids the need to invert an $N \times N$ covariance matrix to compute the posterior.

To choose inducing points that can represent the spread of data in our observations across feature space, we use K-means++ Arthur and Vassilvitskii (2007) with $K = M$ to cluster the feature vectors, then take the cluster centers as inducing points. An alternative approach would be to learn the placement of inducing points as part of the variational inference procedure (?), or by maximizing the variational lower bound on the log marginal likelihood (see next section). However, the former breaks the convergence guarantees, and both approaches may add substantial computational cost. Therefore, in this paper, we show that it is often sufficient to place inducing points up-front, and leaving their optimization for future work.

4.3 SVI for Single-User Preference Learning

To estimate the approximate posterior, we can apply variational inference, which iteratively reduces the KL-divergence between our approximate posterior, $q(s)q(\mathbf{f})q(\mathbf{f}_m)$ and the true posterior, $p(s, \mathbf{f}, \mathbf{f}_m | \mathbf{K}, \alpha_0, \beta_0, \mathbf{y})$, by maximizing a lower bound, \mathcal{L} , on the marginal likelihood, $\log p(\mathbf{y} | \mathbf{K}, \alpha_0, \beta_0)$:

$$\log p(\mathbf{y} | \mathbf{K}, \alpha_0, \beta_0) = \text{KL}(q(\mathbf{f}, \mathbf{f}_m, s) || p(\mathbf{f}, \mathbf{f}_m, s | \mathbf{y}, \mathbf{K}, \alpha_0, \beta_0)) - \mathcal{L}. \quad (21)$$

By taking expectations with respect to the variational q distributions, the lower bound is:

$$\begin{aligned} \mathcal{L} &= \mathbb{E}_{q(\mathbf{f}, \mathbf{f}_m, s)} [\log p(\mathbf{y} | \mathbf{f}) + \log p(\mathbf{f}_m, s | \mathbf{K}, \alpha_0, \beta_0) - \log q(\mathbf{f}_m) - \log q(s)] \\ &= \sum_{p=1}^P \mathbb{E}_{q(\mathbf{f})} [\log p(y_p | f_{a_p}, f_{b_p})] - \frac{1}{2} \left\{ \log |\mathbf{K}_{mm}| - \mathbb{E}[\log s] - \log |\mathbf{S}| - M \right. \\ &\quad \left. + \hat{\mathbf{f}}_m^T \mathbb{E}[s] \mathbf{K}_{mm}^{-1} \hat{\mathbf{f}}_m + \text{Tr}(\mathbb{E}[s] \mathbf{K}_{mm}^{-1} \mathbf{S}) \right\} + \log \Gamma(\alpha) - \log \Gamma(\alpha_0) + \alpha_0 (\log \beta_0) \\ &\quad + (\alpha_0 - \alpha) \mathbb{E}[\log s] + (\beta - \beta_0) \mathbb{E}[s] - \alpha \log \beta, \end{aligned} \quad (22)$$

where $\alpha = \alpha_0 + \frac{M}{2}$ and $\beta = \beta_0 + \frac{\text{Tr}(\mathbf{K}_{mm}^{-1}(\mathbf{S} + \hat{\mathbf{f}}_m \hat{\mathbf{f}}_m^T))}{2}$, and the terms relating to $\mathbb{E}[p(\mathbf{f} | \mathbf{f}_m) - q(\mathbf{f})]$ cancel. The iterative algorithm proceeds by updating each of the q factors in turn, taking expectations with respect to the other factors.

The only term in \mathcal{L} that refers to the observations, \mathbf{y} , is a sum of P terms, each of which refers to one observation only. This means that \mathcal{L} can be maximized iteratively by considering a random subset of observations at each iteration (Hensman et al. 2013). Hence, we replace Equations 19 and 18 for

computing $\hat{\mathbf{f}}_m$ and \mathbf{S} over all observations with a sequence of stochastic updates.

For the i th update, we randomly select observations $\mathbf{y}_i = \{y_p \forall p \in \mathbf{P}_i\}$, where \mathbf{P}_i is random subset of indexes of observations. Rather than using the complete matrices, we perform updates using subsets: \mathbf{Q}_i contains rows and columns for observations in \mathbf{P}_i , \mathbf{K}_{im} and \mathbf{A}_i contain only rows referred to by $y_p \forall p \in \mathbf{P}_i$, \mathbf{G}_i contains rows in \mathbf{P}_i and columns referred to by $a_p \forall p \in \mathbf{P}_i$ and $b_p \forall p \in \mathbf{P}_i$, and $\hat{\mathbf{z}}_i = \{\mathbb{E}[\mathbf{z}_p] \forall p \in \mathbf{P}_i\}$. The update equations optimize the natural parameters of the Gaussian distribution by following the natural gradient (Hensman et al. 2015):

$$\mathbf{S}_i^{-1} = (1 - \rho_i)\mathbf{S}_{i-1}^{-1} + \rho_i \left(\mathbb{E}[s]\mathbf{K}_{mm}^{-1} + w_i\mathbf{K}_{mm}^{-1}\mathbf{K}_{im}^T\mathbf{G}_i^T\mathbf{Q}_i^{-1}\mathbf{G}_i\mathbf{K}_{im}\mathbf{K}_{mm}^{-T} \right) \quad (23)$$

$$\hat{\mathbf{f}}_{m,i} = \mathbf{S}_i \left((1 - \rho_i)\mathbf{S}_{i-1}^{-1}\hat{\mathbf{f}}_{m,i-1} + \rho_i w_i \mathbf{K}_{mm}^{-1} \mathbf{K}_{im}^T \mathbf{G}_i^T \mathbf{Q}_i^{-1} \left(\mathbf{y}_i - \Phi(\mathbb{E}[\mathbf{z}_i]) + \mathbf{G}_i \mathbf{A}_i \hat{\mathbf{f}}_{m,i} \right) \right) \quad (24)$$

where $\rho_i = (i + \text{delay})^{-\text{forgettingRate}}$ is a mixing coefficient that controls the update rate, $w_i = \frac{P}{|\mathbf{P}_i|}$ weights each update according to sample size, and delay and forgettingRate are hyperparameters of the algorithm (Hoffman et al. 2013), .

The scale parameter, s , can also be learned as part of the SVI procedure. Its variational factor, $q(s)$, has the following update equations:

$$\mathbb{E}[s] = \frac{2a_0 + M}{2b} \quad (25)$$

$$\mathbb{E}[\log s] = \Psi(2a_0 + M) - \log(2b), \quad (26)$$

where Ψ is the digamma function.

The complete SVI algorithm is summarized in Algorithm 1. The use of an

Input: Pairwise labels, \mathbf{y} , training item features, \mathbf{x} , test item features \mathbf{x}^*

- 1 Compute kernel matrices \mathbf{K} , \mathbf{K}_{mm} and \mathbf{K}_{nm} given \mathbf{x} Initialise $\mathbb{E}[s]$, $\mathbb{E}[\mathbf{f}]$ and $\hat{\mathbf{f}}_m$ to prior means and \mathbf{S} to prior covariance \mathbf{K}_{mm} ;
- while** \mathcal{L} not converged **do**
- 3 Select random sample, \mathbf{P}_i , of P observations **while** \mathbf{G}_i not converged **do**
- 4 Compute \mathbf{G}_i given $\mathbb{E}[\mathbf{f}_i]$;
- 5 Compute $\hat{\mathbf{f}}_{m,i}$ and \mathbf{S}_i ;
- 6 Compute $\mathbb{E}[\mathbf{f}_i]$;
- end**
- 7 Update $q(s)$ and compute $\mathbb{E}[s]$ and $\mathbb{E}[\log s]$;
- end**
- 8 Compute kernel matrices for test items, \mathbf{K}_{**} and \mathbf{K}_{*m} , given \mathbf{x}^* ;
- 9 Use converged values of $\mathbb{E}[\mathbf{f}]$ and $\hat{\mathbf{f}}_m$ to estimate posterior over \mathbf{f}^* at test points ;

Output: Posterior mean of the test values, $\mathbb{E}[\mathbf{f}^*]$ and covariance, \mathbf{C}^*

Algorithm 1: The SVI algorithm for preference learning with a single user.

inner loop to learn \mathbf{G}_i avoids the need to store the complete matrix, \mathbf{G} . The inferred distribution over the inducing points can be used for predicting the values of test items, $f(\mathbf{x}^*)$:

$$\mathbf{f}^* = \mathbf{K}_{*m} \mathbf{K}_{mm}^{-1} \hat{\mathbf{f}}_m, \quad (27)$$

$$\mathbf{C}^* = \mathbf{K}_{**} + \mathbf{K}_{*m} \mathbf{K}_{mm}^{-1} (\mathbf{S} - \mathbf{K}_{mm}/\mathbb{E}[s]) \mathbf{K}_{*m}^T \mathbf{K}_{mm}^{-1}, \quad (28)$$

where \mathbf{C}^* is the posterior covariance of the test items, \mathbf{K}_{**} is their prior covariance, and \mathbf{K}_{*m} is the covariance between test and inducing points. It is possible to recover the lower bound proposed by Hensman et al. (2015) for classification by generalizing the likelihood to arbitrary nonlinear functions, and omitting terms relating to $p(s|\alpha_0, \beta_0)$ and $q(s)$. However, our approach avoids expensive quadrature methods by linearizing the likelihood to enable analytical updates. We also infer s in a Bayesian manner, rather than treating as a hyper-parameter, which is important for preference learning where s controls the noise level of the observations relative to f .

4.4 SVI for Crowd Preference Learning

We now extend the SVI method to the crowd preference learning model proposed in Section 3.4. To begin with, we extend the variational posterior in Equation 16 to approximate the crowd model defined in Equation 9.

$$p(\mathbf{V}, \mathbf{V}_m, \mathbf{W}, \mathbf{W}_m, \mathbf{t}, \mathbf{t}_m, s_1, \dots, s_C, \sigma | \mathbf{y}, \mathbf{x}, \mathbf{x}_m, \mathbf{u}, \mathbf{u}_m, k, \alpha_0, \beta_0) \approx q(\mathbf{V})q(\mathbf{W})q(\mathbf{t})q(\mathbf{V}_m)q(\mathbf{W}_m)q(\mathbf{t}_m) \prod_{c=1}^C q(s_c)q(\sigma), \quad (29)$$

where \mathbf{u}_m are the feature vectors of inducing points for the users. This approximation factorizes the joint distribution between the latent item factors, \mathbf{V} , the latent user factors, \mathbf{W} , and the common means, \mathbf{t} . The variational factors for the inducing points can be obtained by deriving expectations as follows, beginning with the latent item factors:

$$\begin{aligned} \log q(\mathbf{V}_m) &= \mathbb{E}_{q(\mathbf{W}), q(\mathbf{t})} [\log \mathcal{N}(\mathbf{y}; \tilde{\Phi}(\mathbf{z}), \mathbf{Q})] \\ &\quad + \sum_{c=1}^C \log \mathcal{N}(\mathbf{v}_{m,c}; \mathbf{0}, \mathbf{K}_{v,mm}/\mathbb{E}[s_c]) + \text{const} \\ &= \sum_{c=1}^C \log \mathcal{N}(\mathbf{v}_{m,c}; \hat{\mathbf{v}}_{m,c}, \mathbf{S}_{v,c}). \end{aligned} \quad (30)$$

where the precision is given by:

$$\mathbf{S}_{v,c}^{-1} = \mathbf{K}_{v,mm}^{-1}/\mathbb{E}[s_c] + \mathbf{A}_v^T \mathbf{G}^T \text{diag}(\hat{\mathbf{w}}_{c,\mathbf{u}}^2 + \boldsymbol{\Sigma}_{c,\mathbf{u},\mathbf{u}}) \mathbf{Q}^{-1} \mathbf{G} \mathbf{A}_v, \quad (31)$$

where $\mathbf{A}_v = \mathbf{K}_{v,nm} \mathbf{K}_{v,mm}^{-1}$, $\hat{\mathbf{w}}_c$ and $\boldsymbol{\Sigma}_c$ are the variational mean and covariance of the c th latent user component (defined below in Equations 40 and 39), and the subscript $\mathbf{u} = \{u_p \forall p \in 1, \dots, P\}$ is the vector of user indexes in the observations, \mathbf{y} . The term $\text{diag}(\hat{\mathbf{W}}_{c,j}^2 + \boldsymbol{\Sigma}_{c,j})$ scales the diagonal observation precision, \mathbf{Q}^{-1} , by the latent user factors. We use $\mathbf{S}_{v,c}^{-1}$ to compute the variational means for each row of \mathbf{V}_m as follows:

$$\hat{\mathbf{v}}_{m,c} = \mathbf{S}_{v,c} \mathbf{A}_v^T \mathbf{G}^T \text{diag}(\hat{\mathbf{w}}_{c,j}) \mathbf{Q}^{-1} \left(\mathbf{y} - \Phi(\mathbb{E}[\mathbf{z}]) + \sum_{j=1}^U \mathbf{H}_j (\hat{\mathbf{v}}_c^T \hat{\mathbf{w}}_{c,j}) \right), \quad (32)$$

where $\mathbf{H}_j \in P \times N$ contains partial derivatives of the pairwise likelihood with respect to $F_{i,j} = \hat{v}_{c,i} \hat{w}_{c,j}$, with elements given by:

$$H_{j,p,i} = \Phi(\mathbb{E}[z_p])(1 - \Phi(\mathbb{E}[z_p]))(2y_p - 1)([i = a_p] - [i = b_p])[j = u_p]. \quad (33)$$

This is needed to replace \mathbf{G} in the single-user model, since the vector of latent function values, \mathbf{f} , has been replaced by the matrix \mathbf{F} , where each column of \mathbf{F} corresponds to a single user.

The variational component for the inducing points of the common item mean follows a similar pattern:

$$\begin{aligned} \log q(\mathbf{t}_m) &= \mathbb{E}_{q(\mathbf{V})q(\mathbf{W})}[\log \mathcal{N}(\mathbf{y}; \tilde{\Phi}(\mathbf{z}), \mathbf{Q})] + \mathbb{E}[\log \mathcal{N}(\mathbf{t}_m; \mathbf{0}, \mathbf{K}_{t,mm}/s)] + \text{const} \\ &= \log \mathcal{N}(\mathbf{t}; \hat{\mathbf{t}}, \mathbf{S}_t) \end{aligned} \quad (34)$$

$$\mathbf{S}_t^{-1} = \mathbf{K}_{t,mm}^{-1} / \mathbb{E}[\sigma] + \mathbf{A}_t^T \mathbf{G}^T \mathbf{Q}^{-1} \mathbf{G} \mathbf{A}_t \quad (35)$$

$$\hat{\mathbf{t}}_m = \mathbf{S}_t \mathbf{A}_t^T \mathbf{G}^T \mathbf{Q}^{-1} (\mathbf{y} - \Phi(\mathbb{E}[\mathbf{z}]) + \mathbf{G}(\hat{\mathbf{t}})), \quad (36)$$

where $\mathbf{A}_t = \mathbf{K}_{t,nm} \mathbf{K}_{t,mm}^{-1}$.

***TODO: this is not right – the G term should be the same for all cases.

But it is actually the H term specific above, except we sum over either users or item (these are multiple data points)*** Finally, the latent user factors, \mathbf{W} , require a different linearization matrix, $\mathbf{J} \in P \times U$, containing partial derivatives of the pairwise likelihood with respect to \hat{w}_c . Its elements are given by:

$$J_{p,j} = \Phi(\mathbb{E}[z_p])(1 - \Phi(\mathbb{E}[z_p]))(2y_p - 1)[u_p = j] \quad (37)$$

The variational distribution for the inducing points is then as follows:

$$\begin{aligned} \log q(\mathbf{W}_m) &= \mathbb{E}_{q(\mathbf{V})q(\mathbf{t})}[\log \mathcal{N}(\mathbf{y}; \tilde{\Phi}(\mathbf{z}), \mathbf{Q})] + \sum_{c=1}^C \mathbb{E}[\log \mathcal{N}(\mathbf{w}_c; \mathbf{0}, \mathbf{K}_{w,mm})] + \text{const} \\ &= \sum_{c=1}^C \log \mathcal{N}(\mathbf{w}_c; \hat{\mathbf{w}}_c, \boldsymbol{\Sigma}), \end{aligned} \quad (38)$$

where the variational parameters are:

$$\begin{aligned} \Sigma_c^{-1} = & \mathbf{K}_{w,mm}^{-1} + \mathbf{A}_w^T \left(\mathbf{J}^T \text{diag}(\hat{\mathbf{v}}_{c,a}^2 + \mathbf{S}_{c,a,a} + \hat{\mathbf{v}}_{c,b}^2 + \mathbf{S}_{c,b,b} \right. \\ & \left. - 2\hat{\mathbf{v}}_{c,a}\hat{\mathbf{v}}_{c,b} - 2\mathbf{S}_{c,a,b}) \mathbf{Q}^{-1} \mathbf{J}^T \right) \mathbf{A}_w \end{aligned} \quad (39)$$

$$\begin{aligned} \hat{\mathbf{w}}_{m,c} = & \Sigma_c \mathbf{A}_w^T \left(\mathbf{J}^T \text{diag}(\hat{\mathbf{v}}_{c,a}) - \mathbf{J}^T \text{diag}(\hat{\mathbf{v}}_{c,b}) \right) \mathbf{Q}^{-1} \\ & \left(\mathbf{y} - \Phi(\mathbb{E}[\mathbf{z}]) + \sum_{j=1}^U \mathbf{H}_u(\hat{\mathbf{v}}_c^T \hat{\mathbf{w}}_{c,j}) \right), \end{aligned} \quad (40)$$

where the subscripts $\cdot_a = \{\cdot_{a_p} \forall p \in 1, \dots, P\}$ and $\cdot_b = \{\cdot_{b_p} \forall p \in 1, \dots, P\}$ are lists of indices to the first and second items in the pairs, respectively, and $\mathbf{A}_w = \mathbf{K}_{w,um} \mathbf{K}_{w,mm}^{-1}$.

The equations for the means and covariances can be adapted for stochastic updating by applying weighted sums over the stochastic update and the previous values in the same way as Equation 23 and 24. The stochastic updates for the inducing points of the latent factors depend on expectations with respect to the observed points. As with the single user case, the variational factors at the observed items are independent of the observations given the variational factors of the inducing points (likewise for the observed users):

$$\log q(\mathbf{V}) = \sum_{c=1}^C \log \mathcal{N} \left(\mathbf{v}_c; \mathbf{A}_v \hat{\mathbf{v}}_{m,c}, \frac{\mathbf{K}_v}{\mathbb{E}[s_c]} + \mathbf{A}_v (\mathbf{S}_{m,c} - \frac{\mathbf{K}_{v,mm}}{\mathbb{E}[s_c]}) \mathbf{A}_v \right) \quad (41)$$

$$\log q(\mathbf{t}) = \log \mathcal{N} \left(\mathbf{t}; \mathbf{A}_t \hat{\mathbf{t}}_m, \frac{\mathbf{K}_t}{\mathbb{E}[\sigma]} + \mathbf{A}_t (\mathbf{S}_t - \frac{\mathbf{K}_{t,mm}}{\mathbb{E}[\sigma]}) \mathbf{A}_t \right) \quad (42)$$

$$\log q(\mathbf{W}) = \sum_{c=1}^C \log \mathcal{N}(\mathbf{w}_c; \mathbf{A}_w \hat{\mathbf{w}}_{m,c}, \mathbf{K}_w + \mathbf{A}_w (\Sigma - \mathbf{K}_{w,mm}) \mathbf{A}_w). \quad (43)$$

The expectations for the inverse scales, s_1, \dots, s_c and σ , can be computed using the formulas in Equations 25 and 26 by substituting in the corresponding terms for each \mathbf{v}_c or \mathbf{t} instead of \mathbf{f} . Predictions in the latent component model can be made using Equations 41, 42 and 43 by substituting the covariance terms relating to observation items, \mathbf{x} , and users, \mathbf{u} , with corresponding covariance terms for the prediction items and users.

As with the single user model, the lower bound on the marginal likelihood contains sums over the observations, hence is suitable for stochastic variational

updates:

$$\begin{aligned}
\mathcal{L}_{crowd} = & \sum_{p=1}^P \mathbb{E}_{q(\mathbf{f})} [\log p(y_p | \mathbf{v}_{a_p}^T \mathbf{w}_{a_p} + t_{a_p}, \mathbf{v}_{b_p}^T \mathbf{w}_{b_p} + t_{b_p})] - \frac{1}{2} \left\{ \sum_{c=1}^C \left\{ -M_n - M_u \right. \right. \\
& + \log |\mathbf{K}_{v,mm}| + \log |\mathbf{K}_{w,mm}| - \log |\mathbf{S}_{v,c}| - \mathbb{E}[\log s_c] + \hat{\mathbf{v}}_{m,c}^T \mathbb{E}[s_c] \mathbf{K}_{v,mm}^{-1} \hat{\mathbf{v}}_{m,c} \\
& + \text{Tr}(\mathbb{E}[s_c] \mathbf{K}_{v,mm}^{-1} \mathbf{S}_{v,c}) - \log |\mathbf{S}_c| + \hat{\mathbf{w}}_{m,c}^T \mathbf{K}_{w,mm}^{-1} \hat{\mathbf{w}}_{m,c} + \text{Tr}(\mathbf{K}_{w,mm}^{-1} \mathbf{S}_c) \left. \right\} \\
& - M_n + \log |\mathbf{K}_{t,mm}| - \log |\mathbf{S}_t| - \mathbb{E}[\log \sigma] + \hat{\mathbf{t}}^T \mathbb{E}[\sigma] \mathbf{K}_{t,mm}^{-1} \hat{\mathbf{t}} \\
& + \text{Tr}(\mathbb{E}[\sigma] \mathbf{K}_{t,mm}^{-1} \mathbf{S}_t) \left. \right\} - (C+1)(\log \Gamma(\alpha_0) + \alpha_0(\log \beta_0)) \\
& + \sum_{c=1}^C \left\{ \log \Gamma(\alpha_c) + (\alpha_0 - \alpha_c) \mathbb{E}[\log s_c] + (\beta_c - \beta_0) \mathbb{E}[s_c] - \alpha_c \log \beta_c \right\} \\
& + \log \Gamma(\alpha_\sigma) + (\alpha_0 - \alpha_\sigma) \mathbb{E}[\log \sigma] + (\beta_\sigma - \beta_0) \mathbb{E}[s_c] - \alpha_\sigma \log \beta_\sigma, \quad (44)
\end{aligned}$$

In this section, we proposed an SVI scheme for Bayesian matrix factorization given pairwise observations. The inference scheme can readily be adapted to regression or classification tasks by swapping out the preference likelihood, resulting in different values for \mathbf{G} and \mathbf{H} . We now show how to learn the length-scale parameter required to compute covariances using typical kernel functions, then demonstrate how our method can be applied to learning user preferences or consensus opinion when faced with disagreement.

5 Gradient-based Length-scale Optimization

In the previous sections, we defined preference learning models that incorporate GP priors over the latent functions. The covariances of these GPs are defined by a kernel function k , typically of the following form:

$$k_{\boldsymbol{\theta}}(\mathbf{x}, \mathbf{x}') = \prod_{d=1}^D k_d \left(\frac{|x_d - x'_d|}{l_d}, \boldsymbol{\theta}_d \right) \quad (45)$$

where D is the number of features, l_d is a length-scale hyper-parameter, and $\boldsymbol{\theta}_d$ are additional hyper-parameters for an individual feature kernel, k_d . Each k_d is a function of the distance between the d th feature values in feature vectors \mathbf{x} and \mathbf{x}' . The product over features in k means that data points have high covariance only if the kernel functions, k_d , for all features are high (a soft AND operator). It is possible to replace the product with a sum, causing covariance to increase for every k_d that is similar (a soft OR operator), or other combinations of the individual feature kernels. The choice of combination over features is therefore an additional hyper-parameter.

The length-scale, l_d , controls the smoothness of the function, k_d , across the feature space and the contribution of each feature to the model. If a feature

has a large length-scale, its values, \mathbf{x} , have less effect on $k_{\theta}(\mathbf{x}, \mathbf{x}')$ than if it has a shorter length-scale. Hence, it is important to set l_d to correctly capture feature relevance. A computationally frugal option is the median heuristic:

$$l_{d,MH} = D \text{median}(\{|x_{i,d} - x_{j,d}| \forall i = 1, \dots, N, \forall j = 1, \dots, N\}). \quad (46)$$

The motivation is that the median will normalize the feature, so that features are equally weighted regardless of their scaling. By using a median to perform this normalization, extreme values remain outliers with relatively large distances. Multiplying the median by the number of features, D , prevents the average covariance $k_{\theta}(\mathbf{x}, \mathbf{x}')$ between items from increasing as we add more features using the product kernel in Equation 45. This heuristic has been shown to work reasonably well for the task of comparing distributions (Gretton et al. 2012), but is a simple heuristic with no guarantees of optimality.

An alternative method for setting l_d is Bayesian model selection using the type II maximum likelihood method, which chooses the value of l_d that maximizes the marginal likelihood, $p(\mathbf{y}|\theta)$. Since the marginal likelihoods for our models are intractable, we maximize the value of the variational lower bound, \mathcal{L} , after convergence of the inference algorithm (defined in Equation 22 for a single user, and Equation 44 for the crowd model). Optimizing kernel length-scales in this manner is known as automatic relevance determination (ARD) (Rasmussen and Williams 2006), since the optimal value of l_d depends on the relevance of feature d .

To perform ARD on feature d , we only need to be able to evaluate \mathcal{L} after variational inference has converged with any given value of l_d . However, if we can also compute derivatives of \mathcal{L} with respect to l_d , we can use more efficient gradient-based methods, such as L-BFGS-B (Zhu et al. 1997). These methods perform iterative optimization, using gradients to guide changes for all D length-scales simultaneously. For the single user model, the required gradient with respect to the d th length-scale, l_d , is as follows:

$$\nabla_{l_d} \mathcal{L} = \frac{\partial \mathcal{L}}{\partial \hat{\mathbf{f}}_m} \frac{\partial \hat{\mathbf{f}}_m}{\partial l_d} + \frac{\partial \mathcal{L}}{\partial \mathbf{S}^{-1}} \frac{\partial \mathbf{S}^{-1}}{\partial l_d} + \frac{\partial \mathcal{L}}{\partial a} \frac{\partial a}{\partial l_d} + \frac{\partial \mathcal{L}}{\partial b} \frac{\partial b}{\partial l_d} + \frac{\partial \mathcal{L}}{\partial \mathbf{K}} \frac{\partial \mathbf{K}}{\partial l_d}. \quad (47)$$

We exploit \mathbf{S}^{-1} . The terms involving the variational parameters $\hat{\mathbf{f}}_m$, \mathbf{S} , a and b arise because they depend indirectly on the length-scale through the expectations in the variational factors, $\log q(\cdot)$. However, when the variational inference algorithm has converged, \mathcal{L} is at a maximum, so the partial derivatives of \mathcal{L} with respect to $\hat{\mathbf{f}}_m$, \mathbf{S} , a and b are zero. Hence, after convergence, $\nabla_{l_d} \mathcal{L}$ simplifies to:

$$\nabla_{l_d} \mathcal{L} = \frac{1}{2} \text{tr} \left(\left(\mathbb{E}[\mathbf{s}] (\hat{\mathbf{f}}_m \hat{\mathbf{f}}_m^T + \mathbf{S}^T) \mathbf{K}_{mm}^{-1} - \mathbf{I} \right) \frac{\partial \mathbf{K}_{mm}}{\partial l_d} \mathbf{K}_{mm}^{-1} \right). \quad (48)$$

For the crowd model, we assume that C latent item components, \mathbf{V} have the same kernel function, which is also shared with \mathbf{t} . The gradients with respect

to the length-scale, $l_{w,d}$, for the d th item feature are therefore given by:

$$\begin{aligned} \nabla_{l_{w,d}} \mathcal{L}_{crowd} = & \frac{1}{2} \text{tr} \left(\left(\sum_{c=1}^C \mathbb{E}[s_c] \left\{ \hat{\mathbf{v}}_{m,c} \hat{\mathbf{v}}_{m,c}^T + \mathbf{S}_{v,c}^T \right\} \mathbf{K}_{mm,v}^{-1} - C\mathbf{I} \right) \frac{\partial \mathbf{K}_{mm,v}}{\partial l_{w,d}} \mathbf{K}_{mm,v}^{-1} \right) \\ & + \frac{1}{2} \text{tr} \left(\left(\mathbb{E}[\sigma] (\hat{\mathbf{t}}_m \hat{\mathbf{t}}_m^T + \mathbf{S}_t^T) \mathbf{K}_{mm,t}^{-1} - \mathbf{I} \right) \frac{\partial \mathbf{K}_{mm,t}}{\partial l_{w,d}} \mathbf{K}_{mm,t}^{-1} \right). \end{aligned} \quad (49)$$

The gradients for the d th user feature length-scale, $l_{w,d}$, follows the same form:

$$\nabla_{l_{w,d}} \mathcal{L}_{crowd} = \frac{1}{2} \text{tr} \left(\left(\sum_{c=1}^C \left\{ \hat{\mathbf{w}}_{m,c} \hat{\mathbf{w}}_{m,c}^T + \mathbf{S}_c^T \right\} \mathbf{K}_{mm,w}^{-1} - C\mathbf{I} \right) \frac{\partial \mathbf{K}_{mm,w}}{\partial l_{w,d}} \mathbf{K}_{mm,w}^{-1} \right). \quad (50)$$

The partial derivative of the covariance matrix \mathbf{K}_{mm} with respect to l_d depends on the choice of kernel function. The Matérn $\frac{3}{2}$ function is a widely-applicable, differentiable kernel function that has been shown empirically to outperform other well-established kernels such as the squared exponential, and makes weaker assumptions of smoothness of the latent function (Rasmussen and Williams 2006). It is defined as:

$$k_d \left(\frac{|x_d - x'_d|}{l_d} \right) = \left(1 + \frac{\sqrt{3}|x_d - x'_d|}{l_d} \right) \exp \left(-\frac{\sqrt{3}|x_d - x'_d|}{l_d} \right). \quad (51)$$

Assuming that the kernel functions for each feature, k_d , are combined using a product, as in Equation 45, the partial derivative $\frac{\partial \mathbf{K}_{mm}}{\partial l_d}$ is a matrix, where each entry, i, j , is defined by:

$$\frac{\partial K_{mm,ij}}{\partial l_d} = \prod_{d'=1, d' \neq d}^D k_{d'} \left(\frac{|x_{d'} - x'_{d'}|}{l_{d'}} \right) \frac{3(\mathbf{x}_{i,d} - \mathbf{x}_{j,d})^2}{l_d^3} \exp \left(-\frac{\sqrt{3}|\mathbf{x}_{i,d} - \mathbf{x}_{j,d}|}{l_d} \right), \quad (52)$$

where we assume the use of Equation to combine kernel functions over features using a product

To make use of Equations 48 to 52, we nest the variational algorithm defined in Section 4 inside an iterative gradient-based optimization method. Optimization then begins with an initial guess for all length-scales, l_d , such as the median heuristic. Given the current values of l_d , the optimizer (e.g. L-BFGS-B) runs the VB algorithm to convergence, computes $\nabla_{l_d} \mathcal{L}$, then proposes a new candidate value of l_d . The process repeats until the optimizer converges or reaches a maximum number of iterations, and returns the value of l_d that maximized \mathcal{L} .

Fig. 1 Recovering latent preference functions with increasing levels of noise in different scenarios: (a) single user; (b) ground truth function from crowdsourced labels; (c) latent components from a crowd of users.

6 Experiments

6.1 Methods Compared

We refer to the multi-user variant of our model as *crowd-GPPL*. As baselines, we use GPPL to learn a single preference function from all users' preference labels, (*GPPL-pooled*), and a Gaussian process over the joint feature space of users and items (*GPPL-joint*), as proposed by Guo et al. (2010). For datasets up to 100 users (simulated data, subsamples of the real datasets), we also test separate GPPL instances per user with no collaborative learning (*GPPL-per-user*), but this could not be applied to the real datasets as the computation costs were too high. To test the benefit of using GPs to model item and user features, we also test two further baselines: *crowd-GPPL*\mathbf{u}, which ignores the user features, and *crowd-BMF*, which ignores both user and item features and so does not use GPs at all. For both of these methods, the user covariance matrix, \mathbf{K}_w , in the crowd-GPPL model is replaced by the identity matrix, and for *crowd-BMF*, the item covariance matrices, \mathbf{K}_v and \mathbf{K}_t are also replaced by the identity matrix.

6.2 Simulated Noisy Data

Dataset:

Hypothesis:

7 Comparison of Alternative Methods on Real Data

Dataset: Sushi

Hypothesis:

Setup:

Run 25 repeats of random train/test splits with: [a] 1000 (a la Houlsby 15 training, 5 test pairs per user, $P = 15000, P_{test} = 5000$), Sushi-A (10 items), and [b] 5000 users (a la Khan, 10 training, 1 test pairs per user, $P = 50000, P_{test} = 5000$), Sushi-B (100 items). Evaluate on pairwise labelling error, pairwise label logloss, spearman rank correlation, and runtime.

In all cases, we use the no. pairs per user and no. users to select the subset of data used for training, then test on the remainder.

8 Scalability Experiments

Dataset:

Fig. 2 Average test error for each method on Sushi-A dataset. The above experiments would generate 2 (sushi-A and sushi-B) x 3 (2 x varying no. users + 1 x varying no. pairs) x 3 (two pairwise metrics, one ranking) = 18 plots. This needs trimming! Suggest we resort to: 1 (sushi-A) x 2 (varying no users with two different no. pairs) x 1 = 2 plots. Then we can put any missing results in the text for comparison with other methods, i.e. either logloss or classification error.

Fig. 3 Performance when predicting consensus from crowdsourced argument convincingness judgements: (a) rank correlations; (b) pairwise label classification.

Fig. 4 Performance when predicting personal preferences: pairwise label classification for individual crowdworker convincingness judgements.

Fig. 5 Trade-off between performance gains and number of optimization rounds when using ARD: classification performance when predicting consensus for argument convincingness judgements.

Hypothesis:

List of experiments to include – need new plots for the crowd model:

1. Performance, computation time vs. no. inducing points
2. Computation time vs. dataset size, no. features
3. not done: Performance, computation time, vs update size

9 Performance on Large NLP Dataset

We compare performance of several methods on the Dataset used in Section 8.

Methods:

Hypothesis:

10 Conclusions and Future Work

Acknowledgments

References

- Abbasnejad E, Sanner S, Bonilla EV, Poupart P, et al. (2013) Learning community-based preferences via dirichlet process mixtures of gaussian processes. In: IJCAI, pp 1213–1219
- Adams RP, Dahl GE, Murray I (2010) Incorporating side information in probabilistic matrix factorization with gaussian processes. In: Proceedings of the Twenty-Sixth Conference on Uncertainty in Artificial Intelligence, AUAI Press, pp 1–9
- Ahn S, Korattikara A, Liu N, Rajan S, Welling M (2015) Large-scale distributed bayesian matrix factorization using stochastic gradient mcmc. In: Proceedings of the 21th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, ACM, pp 9–18
- Arthur D, Vassilvitskii S (2007) k-means++: the advantages of careful seeding. In: Proceedings of the Eighteenth Annual ACM-SIAM Symposium on Discrete Algorithms, Society for Industrial and Applied Mathematics, pp 1027–1035
- Bolgár BM, Antal P (2016) Bayesian matrix factorization with non-random missing data using informative gaussian process priors and soft evidences. *Journal of Machine Learning Research* 52:25–36
- Bradley RA, Terry ME (1952) Rank analysis of incomplete block designs: I. the method of paired comparisons. *Biometrika* 39(3/4):324–345
- Busa-Fekete R, Hüllermeier E, El Mesaoudi-Paul A (2018) Preference-based online learning with dueling bandits: A survey. arXiv preprint arXiv:180711398
- Cai C, Sun H, Dong B, Zhang B, Wang T, Wang H (2017) Pairwise ranking aggregation by non-interactive crowdsourcing with budget constraints. In: Distributed Computing Systems (ICDCS), 2017 IEEE 37th International Conference on, IEEE, pp 2567–2568
- Chen G, Zhu F, Heng PA (2018) Large-scale bayesian probabilistic matrix factorization with memo-free distributed variational inference. *ACM Trans Knowl Discov Data* 12(3):1–31:24, DOI 10.1145/3161886, URL <http://doi.acm.org/10.1145/3161886>
- Chen X, Bennett PN, Collins-Thompson K, Horvitz E (2013) Pairwise ranking aggregation in a crowdsourced setting. In: Proceedings of the sixth ACM international conference on Web search and data mining, ACM, pp 193–202
- Chen Y, Suh C (2015) Spectral mle: Top-k rank aggregation from pairwise comparisons. In: International Conference on Machine Learning, pp 371–380
- Chu W, Ghahramani Z (2005) Preference learning with Gaussian processes. In: Proceedings of the 22nd International Conference on Machine learning, ACM, pp 137–144
- Ding N, Qi Y, Xiang R, Molloy I, Li N (2010) Nonparametric bayesian matrix factorization by power-ep. In: Proceedings of the Thirteenth International Conference on Artificial Intelligence and Statistics, pp 169–176
- Fu Y, Hospedales TM, Xiang T, Xiong J, Gong S, Wang Y, Yao Y (2016) Robust subjective visual property prediction from crowdsourced pairwise labels. *IEEE transactions on pattern analysis and machine intelligence* 38(3):563–577
- Gretton A, Sejdinovic D, Strathmann H, Balakrishnan S, Pontil M, Fukumizu K, Sriperumbudur BK (2012) Optimal kernel choice for large-scale two-sample tests. In: Advances in Neural Information Processing Systems, pp 1205–1213
- Guo S, Sanner S, Bonilla EV (2010) Gaussian process preference elicitation. In: Advances in neural information processing systems, pp 262–270
- Hensman J, Fusi N, Lawrence ND (2013) Gaussian processes for big data. In: Proceedings of the Twenty-Ninth Conference on Uncertainty in Artificial Intelligence, AUAI Press, pp 282–290
- Hensman J, Matthews AGdG, Ghahramani Z (2015) Scalable Variational Gaussian Process Classification. In: Proceedings of the Eighteenth International Conference on Artificial Intelligence and Statistics, pp 351–360
- Herbrich R, Minka T, Graepel T (2007) Trueskill: a bayesian skill rating system. In: Advances in neural information processing systems, pp 569–576
- Hoffman MD, Blei DM, Wang C, Paisley JW (2013) Stochastic variational inference. *Journal of Machine Learning Research* 14(1):1303–1347
- Houlsby N, Huszar F, Ghahramani Z, Hernández-Lobato JM (2012) Collaborative Gaussian processes for preference learning. In: Advances in Neural Information Processing

- Systems, pp 2096–2104
- Khan ME, Ko YJ, Seeger MW (2014) Scalable collaborative bayesian preference learning. In: AISTATS, vol 14, pp 475–483
- Kim Y, Kim W, Shim K (2014) Latent ranking analysis using pairwise comparisons. In: Data Mining (ICDM), 2014 IEEE International Conference on, IEEE, pp 869–874
- Kim Y, Kim W, Shim K (2017) Latent ranking analysis using pairwise comparisons in crowdsourcing platforms. *Information Systems* 65:7–21
- Li J, Baba Y, Kashima H (2018) Simultaneous clustering and ranking from pairwise comparisons. In: IJCAI
- Luce RD (1959) On the possible psychophysical laws. *Psychological review* 66(2):81
- Maystre L, Grossglauser M (2017) Just sort it! a simple and effective approach to active preference learning. In: International Conference on Machine Learning, pp 2344–2353
- Minka TP (2001) Expectation propagation for approximate bayesian inference. In: Proceedings of the Seventeenth conference on Uncertainty in artificial intelligence, Morgan Kaufmann Publishers Inc., pp 362–369
- Mosteller F (2006) Remarks on the method of paired comparisons: I. The least squares solution assuming equal standard deviations and equal correlations. In: *Selected Papers of Frederick Mosteller*, Springer, pp 157–162
- Nickisch H, Rasmussen CE (2008) Approximations for binary Gaussian process classification. *Journal of Machine Learning Research* 9(Oct):2035–2078
- Plackett RL (1975) The analysis of permutations. *Applied Statistics* pp 193–202
- Qian L, Gao J, Jagadish H (2015) Learning user preferences by adaptive pairwise comparison. *Proceedings of the VLDB Endowment* 8(11):1322–1333
- Radlinski F, Joachims T (2007) Active exploration for learning rankings from clickthrough data. In: *Proceedings of the 13th ACM SIGKDD international conference on Knowledge discovery and data mining*, ACM, pp 570–579
- Rasmussen CE, Williams CKI (2006) *Gaussian processes for machine learning*. The MIT Press, Cambridge, MA, USA 38:715–719
- Reece S, Roberts S, Nicholson D, Lloyd C (2011) Determining intent using hard/soft data and Gaussian process classifiers. In: *Proceedings of the 14th International Conference on Information Fusion*, IEEE, pp 1–8
- Salakhutdinov R, Mnih A (2008) Bayesian probabilistic matrix factorization using markov chain monte carlo. In: *Proceedings of the 25th international conference on Machine learning*, ACM, pp 880–887
- Shah N, Balakrishnan S, Bradley J, Parekh A, Ramchandran K, Wainwright M (2015) Estimation from pairwise comparisons: Sharp minimax bounds with topology dependence. In: *Artificial Intelligence and Statistics*, pp 856–865
- Shi J, Zheng X, Yang W (2017) Survey on probabilistic models of low-rank matrix factorizations. *Entropy* 19(8):424
- Simpson ED, Gurevych I (2018) Finding convincing arguments using scalable bayesian preference learning. *Transactions of the Association for Computational Linguistics* 6:357–371
- Snelson E, Ghahramani Z (2006) Sparse gaussian processes using pseudo-inputs. In: *Advances in neural information processing systems*, pp 1257–1264
- Steinberg DM, Bonilla EV (2014) Extended and unscented Gaussian processes. In: *Advances in Neural Information Processing Systems*, pp 1251–1259
- Thurstone LL (1927) A law of comparative judgment. *Psychological review* 34(4):273
- Tian Y, Zhu J (2012) Learning from crowds in the presence of schools of thought. In: *Proceedings of the 18th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, ACM, New York, NY, USA, KDD '12, pp 226–234, DOI 10.1145/2339530.2339571, URL <http://doi.acm.org/10.1145/2339530.2339571>
- Uchida S, Yamamoto T, Kato MP, Ohshima H, Tanaka K (2017) Entity ranking by learning and inferring pairwise preferences from user reviews. In: *Asia Information Retrieval Symposium*, Springer, pp 141–153
- Vander Aa T, Chakroun I, Haber T (2017) Distributed bayesian probabilistic matrix factorization. *Procedia Computer Science* 108:1030–1039
- Wang X, Wang J, Jie L, Zhai C, Chang Y (2016) Blind men and the elephant: Thurstonian pairwise preference for ranking in crowdsourcing. In: *Data Mining (ICDM), 2016 IEEE 16th International Conference on, IEEE*, pp 509–518

- Yi J, Jin R, Jain S, Jain A (2013) Inferring Users Preferences from Crowdsourced Pairwise Comparisons: A Matrix Completion Approach. In: First AAAI Conference on Human Computation and Crowdsourcing
- Zhou T, Shan H, Banerjee A, Sapiro G (2012) Kernelized probabilistic matrix factorization: Exploiting graphs and side information. In: Proceedings of the 2012 SIAM international Conference on Data mining, SIAM, pp 403–414
- Zhu C, Byrd RH, Lu P, Nocedal J (1997) Algorithm 778: L-bfgs-b: Fortran subroutines for large-scale bound-constrained optimization. ACM Transactions on Mathematical Software (TOMS) 23(4):550–560