

# Scalable Bayesian Methods for Personalised Preference Learning

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

Received: date

**Abstract** We propose a scalable Bayesian preference learning method for jointly inferring a consensus from crowdsourced pairwise labels and predicting the preferences of individual annotators. Annotators in a crowd may have divergent opinions, making it difficult to identify consensus rankings or ratings from pairwise labels. Limited, noisy data for each user also presents a challenge when predicting the preferences of individuals. We address these challenges by combining matrix factorization with Gaussian processes in a Bayesian approach that accounts for uncertainty arising from sparse data and annotation noise. Previous methods for Gaussian process preference learning (GPPL) do not scale to datasets with large numbers of users, items or pairwise labels, so we propose an inference method using stochastic variational inference (SVI) that can handle constraints on memory and computational costs. Our experiments on a computational argumentation task demonstrate the method’s scalability and show that modeling the preferences of individual annotators in a crowd improves the quality of the inferred consensus. On a benchmark recommendation task, our method is competitive with previous approaches, while the model is able to scale to far larger datasets. We make our software and documentation publicly available for use in future work<sup>1</sup>.

## 1 Introduction

*Preference learning* is the task of learning to compare the values of a set of alternatives according to a particular quality (Fürnkranz and Hüllermeier 2010). The goal may be to rank items by preference, or given a pair of items, to select the item with the highest utility. For many preference learning tasks, annotators have divergent opinions on the correct rankings or utilities, which impedes the acquisition of training data. For example, in the field of argument mining, one goal is to identify convincing arguments from a corpus of documents (Habernal and Gurevych 2016). Whether a particular argument is convincing or not depends

---

Address(es) of author(s) should be given

<sup>1</sup> <https://github.com/UKPLab/tac12018-preference-convincing/tree/crowdGPPL>

on the reader’s point of view and prior knowledge (Lukin et al. 2017). Similarly, recommender systems can perform better if they make recommendations tailored to a specific user (Resnick and Varian 1997). Crowdsourcing is frequently used as a cost-effective source of labelled data, yet disagreements between annotators must be resolved to obtain a gold-standard training set, typically requiring redundant labelling and increased annotation costs (Snow et al. 2008; Banerji et al. 2010; Gaunt et al. 2016). Therefore, we require solutions for predicting individual preferences or producing a gold standard from crowdsourced data that take into account the differences of opinion between annotators and users.

A preference model can be trained using numerical ratings, yet each annotator may interpret the values differently and may label inconsistently depending on the order in which they annotate items (Ovadia 2004; Yannakakis and Hallam 2011): a score of -1, say, from one annotator may be equivalent to a -20 from another. An alternative is *pairwise labelling*, in which the annotator compares pairs of items and selects the preferred one in each pair. Making pairwise choices places lower cognitive load on annotators than numerical ratings (Yang and Chen 2011), and facilitates the total sorting of items, since it avoids cases where two items have the same value (Kendall 1948; Kingsley and Brown 2010). Besides explicit annotations, pairwise preference labels can be extracted from user behaviour logs, such as when a user selects one item from a list in preference to others (Joachims 2002). However, these implicit annotations can be noisy, as are crowdsourced pairwise labels (Habernal and Gurevych 2016). Therefore, while pairwise labels provide a valuable form of training data, preference learning methods must be able to aggregate noisy sources of pairwise data, such as that obtained from crowds. We henceforth refer to the annotators, users or implicit data sources simply as *users* whose preferences we wish to model. Likewise, the term *items* refers to any type of instance that users may express preferences over, and could be states or actions as well as physical objects.

In recommender systems, information from different users is aggregated using *collaborative filtering*, which predicts a user’s preferences for an item they have not annotated using the observed preferences of similar users for that item (Resnick and Varian 1997). A typical approach is to represent observed ratings in a user-item matrix, then apply *matrix factorization* (Koren et al. 2009) to decompose a user-item matrix into two low-dimensional matrices. Users and items with similar observed ratings have similar row vectors in the low-dimensional matrices. Multiplying the low-dimensional representations predicts ratings for unseen user-item pairs. However, traditional approaches are not able to handle pairwise labels, nor extrapolate effectively to new users or items.

The limited amount of noisy data for each user, as well as the desire to aggregate data from multiple users, motivates a Bayesian approach that can account for uncertainty in the model. Matrix factorization, has been shown to benefit from a Bayesian treatment when data is sparse (Salakhutdinov and Mnih 2008). Previous work has also introduced a Bayesian approach for combining crowdsourced preferences, but this models only ground truth rather than individual preferences and cannot make predictions for new users or items (Chen et al. 2013). For crowdsourced classification tasks, Simpson et al. (2017) show that modelling item features using a Gaussian process can improve performance, but this has not yet been adapted for preference learning. Gaussian processes have also been used in previous work to provide a Bayesian framework for preference learning (Chu and

Ghahramani 2005; Houlby et al. 2012; Khan et al. 2014), but these methods do not scale to large numbers of items, users, or pairwise labels as it is unclear how to parallelize them or constrain their memory requirements.

In this paper, we propose a scalable Bayesian approach to pairwise preference learning with crowds, whose members may be annotators or users. Our approach is designed for preference learning over items, rather than over labels (see discussion in Fürnkranz and Hüllermeier (2010)), and jointly models personal preferences and the consensus of a crowd by combining matrix factorization and Gaussian processes. To enable inference at the scale of large, real-world datasets, we derive a stochastic variational inference scheme (Hoffman et al. 2013) for our approach. By providing a scalable Bayesian approach we open preference learning up to novel applications for which annotations are sparse, noisy and biased, and where the number of users, items and pairwise labels is large. Our empirical evaluation demonstrates the scalability of our approach, and its ability to predict both personal preferences and an objective gold-standard given crowdsourced data. We also improve performance over the previous state-of-the-art (Simpson and Gurevych 2018) on a crowdsourced argumentation dataset.

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

## 2 Related Work

### 2.1 Aggregating Pairwise Labels from a Crowd

To obtain a ranking from pairwise labels, many preference learning methods model the user’s choices as a random function of the latent utility of the items (Thurstone 1927). An example is the method of Herbrich et al. (2007), which learns the skill of chess players from match outcomes by treating them as noisy pairwise labels. Recent work on this type of approach has analyzed bounds on error rates (Chen and Suh 2015), sample complexity (Shah et al. 2015), and joint models for ranking and clustering from pairwise comparisons (Li et al. 2018).

The problem of disagreement between annotators in a crowd was addressed by Chen et al. (2013) and Wang et al. (2016), using Bayesian approaches that learn the individual accuracy of each worker. However, they do not exploit item features to mitigate data sparsity. In contrast, Gaussian processes preference learning (*GPPL*) uses item features to make predictions for unseen items and share information between similar items (Chu and Ghahramani 2005). GPPL was used by Simpson and Gurevych (2018) to aggregate crowdsourced pairwise labels, but assumes the same level of noise for all annotators and ignores the effect of divergent opinions. To crowdsource sentiment annotations, Kiritchenko and Mohammad (2016) propose to use an extension of pairwise comparison known as *best-worst scaling*, in which the annotator selects best and worst items from a set. They apply a simple counting technique to infer a ranking over the items, which requires each item to have a sufficient number of comparisons.

How do we make it clear that this is not the case? – This paper extended the pairwise ranking model in Xi Chen’s work, Pairwise Ranking Aggregation in a Crowd-sourced Setting with Gaussian Process. The preference learning

A popular method for predicting pairwise labels of new items given their features is *SVM-rank* (Joachims 2002). For crowdsourced data, Fu et al. (2016) show that performance is improved by identifying outliers in crowdsourced data that correspond to probable errors. Uchida et al. (2017) extend SVM-rank to account for different levels of confidence in each pairwise annotation expressed by the annotators. However, besides modeling labeling noise, the previous work on crowdsourced preference learning does not account for the effect of divergence of opinion on the inferred preferences.

## 2.2 Bayesian Methods for Inferring Individual Preferences

[https://github.com/traagvictor/Ascent/blob/master/missing\\_paper?](https://github.com/traagvictor/Ascent/blob/master/missing_paper?)

Section 2.2: Define or give an example of 'input features' when introduced.

Ascent, a tool for aggregating preferences from a crowd to identify a consensus, we also wish to predict the preferences of individual users. Yi et al. (2013) and Kim et al. (2014) address this task by learning multiple latent rankings and inferring the preference of each user toward those rankings, while Salimans et al. (2012) use Bayesian matrix factorization to identify multiple latent rankings. However, none of these approaches exploit item features to remedy labeling errors or generalize to new test items. In contrast, Guo et al. (2010) propose a joint Gaussian process over the space of users and features. Since this scales cubically in the number of users, Abbasnejad et al. (2013) propose to cluster the users into behavioural groups. However, distinct clusters do not allow for collaborative learning between users with partially overlapping preferences, e.g. two users may both like one genre of music, while having different preferences over other genres. Khan et al. (2014) instead learn a GP for each user, and combine them with matrix factorization to perform collaborative filtering. However, this approach does not model the relationship between input features and the latent factors, does not place a prior over item factors, and does not scale to very large sets of users. An alternative is *Collaborative GPPL* (Houlsby et al. 2012), which uses a latent factor model, where each latent factor has a Gaussian process prior. This allows the model to take advantage of the input features of users and items when learning the latent factors. Each individual's preferences are then represented by a mixture of latent functions. Using matrix factorization in combination with GP priors is therefore an effective way to model the individual preferences of users while making use of input features to generalize to test cases. However, none of these approaches considers a consensus preference function, and more scalable inference is needed for datasets containing thousands of items or users.

## 2.3 Scalable Approximate Bayesian Inference

Another claim from the authors is the SVI for large-scale crowdGPPL, which is due to limitation of GP. Indeed, the online strategy has been well-studied in crowdBT, crowdPL, the

Recent work on scalable Bayesian matrix factorization focuses on distributing and parallelizing inference but is not directly applicable to Gaussian processes (Ahn et al. 2015; Vander Aa et al. 2017; Chen et al. 2018). Models that combine Gaussian processes with non-Gaussian likelihoods require approximate inference methods that often scale poorly with the amount of training data available. Established methods such as the Laplace approximation and expectation propagation (Rasmussen and Williams 2006) have computational complexity  $\mathcal{O}(N^3)$  with  $N$  data points and memory complexity  $\mathcal{O}(N^2)$ . For collaborative GPPL, Houlsby et al.

(2012) propose a kernel for pairwise preference learning and use a sparse *generalized fully independent training conditional* (GFITC) approximation (Snelson and Ghahramani 2006) to reduce the computational complexity to  $\mathcal{O}(PM^2 + UM^2)$  and memory complexity to  $\mathcal{O}(PM + UM)$ , where  $P$  is the number of pairwise labels,  $M \ll P$  is a fixed number of inducing points, and  $U$  is the number of users. However, this is not sufficiently scalable for very large numbers of items, users or pairs, as there is no clear way to parallelize it or to limit memory consumption. The GP over pairs also makes it difficult to extract posteriors for latent function values for individual items, and prevents mixing pairwise training labels with observed ratings in future data aggregation settings.

To handle large numbers of pairwise labels, Khan et al. (2014) develop a variational EM algorithm and sub-sample pairwise data rather than learning from the complete training set. An alternative is *Stochastic variational inference (SVI)* (Hoffman et al. 2013), which updates an approximation using a different random sample at each iteration. This allows the approximation to make use of all training data over a number of iterations, while limiting training costs per iteration. SVI has been successfully applied to Gaussian process regression (Hensman et al. 2013) and classification (Hensman et al. 2015), and provides a convenient framework for sparse approximation. An SVI method was also developed for preference learning, which places a GP over items rather than pairs, but does not model multiple users' preferences (Simpson and Gurevych 2018). This paper provides the first full derivation of this approach, as well as providing the first application of SVI to Bayesian matrix factorization as part of a new model that accounts for individual user preferences.

add in  
table of  
notations  
(maybe  
in the  
ap-  
pendix)

### 3 Bayesian Preference Learning for Crowds

For a pair of items,  $a$  and  $b$ , we write  $a \succ b$  to symbolize that  $a$  is preferred to  $b$ . A pairwise label has value  $y(a, b) = 1$  if  $a \succ b$  and 0 otherwise. Thurstone (1927) proposed that the likelihood of pairwise label  $y(a, b) = 1$  increases as the difference between the utilities of  $a$  and  $b$  increases. The uncertainty in the likelihood  $p(y(a, b) = 1)$  accommodates labeling errors as well as variability in a user's judgments. Here, we assume the utility to be a function of the items' features,  $f(\mathbf{x}_a) : \mathbb{R}^D \mapsto \mathbb{R}$ , where  $\mathbf{x}_a$  is a vector representation of the features of item  $a$  of length  $D$ . Pairwise labels can be modeled using a logistic likelihood defined by the Bradley-Terry model (Bradley and Terry 1952; Plackett 1975; Luce 1959), or a probit likelihood given by the Thurstone-Mosteller model, also known as *Thurstone case V* (Thurstone 1927; Mosteller 2006). We use the latter for our model as it enables  $f$  to be marginalized analytically to compute a posterior  $p(y(a, b) = 1 | \mathbf{y})$ , where  $\mathbf{y}$  is a set of pairwise training labels. Noise in the observations is explained by a Gaussian-distributed noise term,  $\delta \sim \mathcal{N}(0, 0.5)$ :

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

Integrating out the unknown values of  $\delta_a$  and  $\delta_b$  gives:

$$\begin{aligned} & p(y(a, b) | f(\mathbf{x}_a), f(\mathbf{x}_b)) \\ &= \iint p(y(a, b) | f(\mathbf{x}_a) + \delta_a, f(\mathbf{x}_b) + \delta_b) \mathcal{N}(\delta_a; 0, 0.5) \mathcal{N}(\delta_b; 0, 0.5) d\delta_a d\delta_b \\ &= \Phi(f(\mathbf{x}_a) - f(\mathbf{x}_b)) = \Phi(z), \end{aligned} \quad (2)$$

where  $z = f(\mathbf{x}_a) - f(\mathbf{x}_b)$ , and  $\Phi$  is the cumulative distribution function of the standard normal distribution. This likelihood is the basis of Gaussian process preference learning (GPPL) (Chu and Ghahramani 2005). Our formulation differs in that instead of learning the variance of  $\delta$ , we fix it to 0.5 and scale  $f$  to vary the uncertainty in the pairwise labels, which conveniently leads to  $z$  having a denominator of 1. If we have Gaussian distributions over  $f(\mathbf{x}_a)$  and  $f(\mathbf{x}_b)$ , they can be marginalized in a simple manner by modifying  $z$ :

$$\hat{z} = \frac{\hat{f}_a - \hat{f}_b}{\sqrt{1 + \sigma_a^2 + \sigma_b^2 - 2\sigma_{a,b}^2}} \quad (3)$$

where  $\hat{f}_a$  and  $\hat{f}_b$  are the expected values of  $f(\mathbf{x}_a)$  and  $f(\mathbf{x}_b)$  respectively,  $\sigma_a^2$  and  $\sigma_b^2$  are the corresponding variances, and  $\sigma_{a,b}^2$  is the covariance between  $f(\mathbf{x}_a)$  and  $f(\mathbf{x}_b)$ . The next section introduces a model for the utility function,  $f$ .

### 3.1 Single User Preference Learning

We assume a Gaussian process prior over the utility function,  $f \sim \mathcal{GP}(0, k_\theta/s)$ , where  $k_\theta$  is a kernel function with hyper-parameters  $\theta$ , and  $s$  is an inverse scale drawn from a gamma prior,  $s \sim \mathcal{G}(\alpha_0, \beta_0)$ , with shape  $\alpha_0$  and scale  $\beta_0$ . 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 2). The kernel function takes numerical item features as inputs and determines the covariance between values of  $f$  for different items. The choice of kernel function and its hyperparameters controls the shape and smoothness of the function across the feature space and is often treated as a model selection problem. Typical choices of kernel function include the *squared exponential* or *Matérn* functions (Rasmussen and Williams 2006), which produce higher covariance between items with similar feature values. These kernel functions make minimal assumptions and are effective in a wide range of tasks. We use the kernel function  $k_\theta$  to compute the covariance matrix  $K_\theta$ , between a set of  $N$  observed items with features  $\{\mathbf{x}_1, \dots, \mathbf{x}_N\}$ .

Given a set of  $P$  pairwise labels,  $\mathbf{y} = \{y_1, \dots, y_P\}$ , where  $y_p = y(a_p, b_p)$  is the preference label for items  $a_p$  and  $b_p$ , we can write the joint distribution over all variables as follows:

$$p(\mathbf{y}, \mathbf{f}, s | k_\theta, \alpha_0, \beta_0) = \prod_{p=1}^P p(y_p | \mathbf{f}) \mathcal{N}(\mathbf{f}; \mathbf{0}, \mathbf{K}_\theta/s) \mathcal{G}(s; \alpha_0, \beta_0) \quad (4)$$

where  $\mathbf{f} = \{f(\mathbf{x}_1), \dots, f(\mathbf{x}_N)\}$  are the utilities of the  $N$  items referred to by  $\mathbf{y}$ , and  $p(y_p | \mathbf{f}) = \Phi(z_p)$  is the pairwise likelihood (Equation 2). We henceforth refer to this single-user model as *GPPL*.

### 3.2 Crowd Preference Learning

When there are multiple users, jointly modeling their different preferences may improve consensus predictions and improve our predictions for each user when faced with sparse data. We represent utilities in a matrix,  $\mathbf{F} \in \mathbb{R}^{N \times U}$ , where  $N$  rows correspond to items and  $U$  columns correspond to users. If we factorize this matrix, we obtain two low-dimensional matrices, one for users,  $\mathbf{W} \in \mathbb{R}^{C \times U}$ , and one for the items,  $\mathbf{V} \in \mathbb{R}^{C \times N}$ , where  $C$  is the number of latent components:  $\mathbf{F} = \mathbf{V}^T \mathbf{W}$ . The column  $\mathbf{v}_{\cdot,a}$  of  $\mathbf{V}$ , and the column  $\mathbf{w}_{\cdot,j}$  of  $\mathbf{W}$ , are latent vector representations of item  $a$  and a user  $j$ , respectively. Latent components correspond to utility functions for certain items shared by multiple users. These could represent, for example, in the case of book recommendation, interests in a particular genre of book. For the multi-user case, we assume that there are  $C$  latent functions of item features,  $v_c \sim \mathcal{GP}(\mathbf{0}, k_\theta/s_c)$ , and  $C$  latent functions of user features,  $w_c \sim \mathcal{GP}(\mathbf{0}, k_\eta)$ , where  $k$  is a kernel function,  $s_c$  is an inverse scale for the item features,  $\theta$  are the kernel hyperparameters for the items and  $\eta$  are the kernel hyperparameters for the users. The matrices  $\mathbf{V}$  and  $\mathbf{W}$  are evaluations of these functions at the points corresponding to the users and items observed during training. Since our goal is to infer a consensus from a crowd as well as to model individual users' preferences, we introduce a consensus utility function over item features,  $t \sim \mathcal{GP}(\mathbf{0}, k_\theta/s_t)$ , that is shared across all users, with values  $\mathbf{t} = \{t(\mathbf{x}_1), \dots, t(\mathbf{x}_N)\}$  for the training items. Therefore, the latent preference function,  $f$ , is a weighted sum over latent functions:

$$f(\mathbf{x}_a, \mathbf{u}_j) = \sum_{c=1}^C v_c(\mathbf{x}_a) w_c(\mathbf{u}_j) + t(\mathbf{x}_a), \quad (5)$$

where  $\mathbf{u}_j$  are the features of user  $j$  and  $\mathbf{x}_a$  are the features of item  $a$ . For simplicity, we assume fixed values of  $C$  in this paper. The Bayesian approach avoids overfitting by inferring  $s_c \approx 0$  with high probability for any dimensions that are not required to model the data.

We combine the matrix factorization method with the preference likelihood of Equation 2 to obtain a joint preference model for multiple users, *crowdGPPL*:

$$p(\mathbf{y}, \mathbf{V}, \mathbf{W}, \mathbf{t}, s_1, \dots, s_C, s_t | k_\theta, k_\eta, \alpha_0, \beta_0) = \prod_{p=1}^P \Phi(z_p) \mathcal{N}(\mathbf{t}; \mathbf{0}, \mathbf{K}_\theta/s_t) \\ \mathcal{G}(s_t; \alpha_0, \beta_0) \prod_{c=1}^C \{\mathcal{N}(\mathbf{v}_c; \mathbf{0}, \mathbf{K}_\theta/s_c) \mathcal{N}(\mathbf{w}_c; \mathbf{0}, \mathbf{L}_\eta) \mathcal{G}(s_c; \alpha_0, \beta_0)\}, \quad (6)$$

where  $z_p = \mathbf{v}_{\cdot, a_p}^T \mathbf{w}_{\cdot, u_p} + t_{a_p} - \mathbf{v}_{\cdot, b_p}^T \mathbf{w}_{\cdot, u_p} - t_{b_p}$ ,  $s_t$  is the inverse scale of  $t$ , index  $p$  refers to a tuple,  $\{u_p, a_p, b_p\}$ , which identifies the user and a pair of items, and  $\mathbf{L}_\eta$  is the prior covariance between user feature vectors computed using the kernel function  $k_\eta$ .

## 4 Scalable Inference

In the single user case, the goal is to infer the posterior distribution over the utilities of test items,  $\mathbf{f}^*$ , given a set of pairwise training labels,  $\mathbf{y}$ . In the multi-user

The only new point is equation (5), but it is not clear why introducing  $t$ . I expected the authors should have more discussion on this.

Page 7, equation (5): you might want to cite precedence to this form of equation e.g. K. Mo, E. Zhong, and Q. Yang. Cross-task crowd-sourcing. In Proc.

case, we aim to find the posterior over the matrix  $\mathbf{F}^* = \mathbf{V}^{*T} \mathbf{W}^*$  of utilities for test items and test users. The non-Gaussian likelihood makes exact inference intractable, hence previous work has used the Laplace approximation for the single user case (Chu and Ghahramani 2005) or 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 classification (Nickisch and Rasmussen 2008). EP and VB approximate the true posterior with a simpler, factorized distribution that can be learned using an iterative algorithm. For crowdGPPL, the true posterior is multi-modal, since the latent factors can be re-ordered arbitrarily without affecting  $\mathbf{F}$ , causing a *non-identifiability problem*. EP would average these modes and produce uninformative predictions over  $\mathbf{F}$ , so Houlsby et al. (2012) incorporate a VB step that approximates a single mode. A drawback of EP is that unlike VB, convergence is not guaranteed (Minka 2001).

Exact inference for a Gaussian process has computational complexity  $\mathcal{O}(N^3)$  and memory complexity  $\mathcal{O}(N^2)$ , where  $N$  is the number of data points. 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 cut to  $\mathcal{O}(NM^2)$ , and the memory complexity to  $\mathcal{O}(NM)$ . 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. The sparse approximation used by Houlsby et al. (2012) for the collaborative GP is the *generalized fully independent training conditional* (GFITC) (Snelson and Ghahramani 2006). In practice, GFITC is unsuitable for datasets with more than a few thousands points, as the computational and memory costs cannot be constrained and GFITC is not amenable to distributed computation (Hensman et al. 2015). We derive a more scalable approach for GPPL and crowdGPPL using stochastic variational inference (SVI), an iterative scheme that limits the computational and memory costs at each iteration (Hoffman et al. 2013). First, we define an approximate likelihood that enables the SVI method.

#### 4.1 Approximate Pairwise Likelihood

To obtain a tractable approximate posterior, we begin by approximating the expected preference likelihood (Equation 2) with a Gaussian:

$$p(\mathbf{y}|\mathbf{f}) = \mathbb{E} \left[ \prod_{p=1}^P \Phi(z_p) \right] = \prod_{p=1}^P \Phi(\hat{z}_p) \approx \mathcal{N}(\mathbf{y}; \Phi(\hat{\mathbf{z}}), \mathbf{Q}), \quad (7)$$

where  $\mathbf{Q}$  is a diagonal noise covariance matrix and  $\hat{\mathbf{z}} = \{\hat{z}_1, \dots, \hat{z}_P\}$ , and  $\hat{z}_p$  is defined by Equation 3. Since  $\Phi(\hat{z}_p)$  defines a bernoulli distribution, for which the conjugate prior is a beta distribution, we moment match the diagonal entries of  $\mathbf{Q}$  to the expected variance of a bernoulli distribution as follows:

$$Q_{p,p} = \mathbb{E}_f[\Phi(z_p)(1 - \Phi(z_p))] \approx \frac{(y_p + \gamma)(1 - y_p + \lambda)}{(\gamma + \lambda + 1)} - \frac{(y_p + \gamma)(1 - y_p + \lambda)}{(\gamma + \lambda + 1)^2(\gamma + \lambda + 2)}, \quad (8)$$

Page 8, line 12: VB is only guaranteed to converge for conjugate distributions.

Page 8, second paragraph: you criticise GFITC for not being decentralisable. It is not clear to me that crowdGPPL is decentralisable.

got really confused by Section 4.1. In equation (7),  $p(\mathbf{y}|\mathbf{f})$  is a product of the cdfs. Where does the expectation come from just after the first equal sign? The



where  $\gamma$  and  $\lambda$  are parameters of a beta distribution with the same variance as the prior,  $p(\Phi(z_p)|\mathbf{K}_\theta, \alpha_0, \beta_0)$ , and are estimated using numerical integration.

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 expectation of  $\mathbf{f}$  for GPPL (for crowdGPPL, replace  $\mathbf{f}$  with  $\mathbf{F}$  in the following):

$$\Phi(\mathbf{z}) \approx \tilde{\Phi}(\mathbf{z}) = \mathbf{G}(\mathbf{f} - \mathbb{E}[\mathbf{f}]) + \Phi(\hat{\mathbf{z}}), \quad (9)$$

$$G_{p,i} = \Phi(\mathbb{E}[z_p])(1 - \Phi(\hat{z}_p))(2y_p - 1)([i = a_p] - [i = b_p]) \quad (10)$$

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 circular dependency between the posterior mean of  $\mathbf{f}$  and  $\mathbf{G}$ , which can be estimated iteratively using variational inference (Steinberg and Bonilla 2014), as we describe below. Linearization makes the approximate likelihood conjugate to the prior,  $\mathcal{N}(\mathbf{f}; \mathbf{0}, \mathbf{K}_\theta/s)$ , so that the approximate posterior over  $\mathbf{f}$  is also Gaussian. The Gaussian likelihood approximation and linearization also appear in GP inference methods based on expectation propagation (Rasmussen and Williams 2006) and the extended Kalman filter (Reece et al. 2011; Steinberg and Bonilla 2014). We can now use the approximate likelihood to derive equations for SVI.

#### 4.2 SVI for Single User GPPL

We introduce a sparse approximation to the Gaussian process that allows us to limit the size of the covariance matrices we need to work with. To do this, we introduce a set of  $M \ll N$  inducing items with inputs  $\mathbf{x}_m$ , utilities  $\mathbf{f}_m$ , covariance  $\mathbf{K}_{mm}$ , and covariance between the observed and inducing items,  $\mathbf{K}_{nm}$ . For clarity, we omit  $\theta$  from this point on. We assume an approximation to the posterior over the inducing and training items that factorises between different sets of latent variables:

$$p(\mathbf{f}, \mathbf{f}_m, s | \mathbf{y}, \mathbf{x}, \mathbf{x}_m, k_\theta, \alpha_0, \beta_0) \approx q(\mathbf{f}, \mathbf{f}_m, s) = q(s)q(\mathbf{f})q(\mathbf{f}_m), \quad (11)$$

where  $q(\cdot)$  are *variational factors*, which we define below. Each factor for a subset of latent variables,  $\mathbf{z}$ , takes the form  $\ln q(\mathbf{z}) = \mathbb{E}_{\mathbf{z}}[\ln p(\mathbf{z}, \mathbf{x}, \mathbf{y})]$ , that is, the expectation with respect to all other latent variables,  $\mathbf{z}/\mathbf{z}$ , of the log joint distribution of the observations and the current subset of latent variables,  $\mathbf{z}$ .

We marginalize  $\mathbf{f}$  to obtain the factor for  $\mathbf{f}_m$ :

$$\begin{aligned} \ln q(\mathbf{f}_m) &= \ln \mathcal{N}(\mathbf{y}; \tilde{\Phi}(\mathbf{z}), \mathbf{Q}) + \ln \mathcal{N}(\mathbf{f}_m; \mathbf{0}, \mathbf{K}_{mm}/\mathbb{E}[s]) + \text{const}, \\ &= \ln \mathcal{N}(\mathbf{f}_m; \hat{\mathbf{f}}_m, \mathbf{S}), \end{aligned} \quad (12)$$

where the variational parameters  $\hat{\mathbf{f}}_m$  and  $\mathbf{S}$  are computed using the iterative SVI procedure described below. We choose an approximation of  $q(\mathbf{f})$  that depends only on the inducing point utilities,  $\mathbf{f}_m$ , and is independent of the observations:

$$\ln q(\mathbf{f}) = \ln \mathcal{N}(\mathbf{f}; \mathbf{A}\hat{\mathbf{f}}_m, \mathbf{K} + \mathbf{A}(\mathbf{S} - \mathbf{K}_{mm}/\mathbb{E}[s])\mathbf{A}^T), \quad (13)$$

Page 9, Section 4.2: Please state up front that you are proposing a mean-field variational approximation in equation (11). This will then justify why  $\mathbb{E}[s]$  appears in equation (13).

where  $\mathbf{A} = \mathbf{K}_{nm}\mathbf{K}_{mm}^{-1}$ . This means we no longer need to invert an  $N \times N$  covariance matrix to compute  $q(\mathbf{f})$ . The factor  $q(s)$  is also modified to depend on the inducing points:

$$\ln q(s) = \mathbb{E}_{q(\mathbf{f}_m)}[\ln \mathcal{N}(\mathbf{f}_m | \mathbf{0}, \mathbf{K}_{mm}/s)] + \ln \mathcal{G}(s; \alpha_0, \beta_0) + \text{const} = \ln \mathcal{G}(s; \alpha, \beta), \quad (14)$$

where  $\alpha = \alpha_0 + \frac{M}{2}$  and  $\beta = \beta_0 + \frac{\text{tr}(\mathbf{K}_{mm}^{-1}(S + \hat{\mathbf{f}}_m \hat{\mathbf{f}}_m^T))}{2}$ .

The location of inducing points can be learned as part of the variational inference procedure (Hensman et al. 2015), or by optimizing a bound on the log marginal likelihood. However, the former breaks the convergence guarantees, and both approaches may add substantial computational cost. We find that we are able to obtain good performance by choosing inducing points up-front using K-means++ (Arthur and Vassilvitskii 2007) with  $K = M$  to cluster the feature vectors, then taking the cluster centers as inducing points that represent the spread of observations across feature space.

We can apply variational inference to iteratively reduce the KL-divergence between our approximate posterior and the true posterior (both stated in Equation 11) by maximizing a lower bound,  $\mathcal{L}$ , on the log marginal likelihood (see also the detailed equations in Appendix A):

$$\begin{aligned} \ln p(\mathbf{y} | \mathbf{K}, \alpha_0, \beta_0) &= \text{KL}(q(\mathbf{f}, \mathbf{f}_m, s) || p(\mathbf{f}, \mathbf{f}_m, s | \mathbf{y}, \mathbf{K}, \alpha_0, \beta_0)) + \mathcal{L} \\ \mathcal{L} &= \mathbb{E}_{q(\mathbf{f})}[\ln p(\mathbf{y} | \mathbf{f})] + \mathbb{E}_{q(\mathbf{f}_m, s)}[\ln p(\mathbf{f}_m, s | \mathbf{K}, \alpha_0, \beta_0) - \ln q(\mathbf{f}_m) - \ln q(s)]. \end{aligned} \quad (15)$$

We optimize  $\mathcal{L}$  by initializing the  $q$  factors randomly, then updating each one 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). For the  $i$ th update of  $q(\mathbf{f}_m)$ , 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. We then perform updates using  $\mathbf{Q}_i$  (rows and columns of  $\mathbf{Q}$  for observations in  $\mathbf{P}_i$ ),  $\mathbf{K}_{im}$  and  $\mathbf{A}_i$ , (rows of  $\mathbf{K}_{nm}$  and  $\mathbf{A}$  referred to by  $\{y_p \forall p \in \mathbf{P}_i\}$ ),  $\mathbf{G}_i$  (rows of  $\mathbf{G}$  in  $\mathbf{P}_i$  and columns referred to by any items  $\{a_p \forall p \in \mathbf{P}_i\} \cup \{b_p \forall p \in \mathbf{P}_i\}$ ), 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):

$$\mathbf{S}_i^{-1} = (1 - \rho_i) \mathbf{S}_{i-1}^{-1} + \rho_i \left( \mathbb{E}[s] \mathbf{K}_{mm}^{-1} + \pi_i \mathbf{A}_i^T \mathbf{G}_i^T \mathbf{Q}_i^{-1} \mathbf{G}_i \mathbf{A}_i \right) \quad (16)$$

$$\hat{\mathbf{f}}_{m,i} = \mathbf{S}_i \left( (1 - \rho_i) \mathbf{S}_{i-1}^{-1} \hat{\mathbf{f}}_{m,i-1} + \rho_i \pi_i \mathbf{A}_i^T \mathbf{G}_i^T \mathbf{Q}_i^{-1} \left( \mathbf{y}_i - \Phi(\mathbb{E}[\mathbf{z}_i]) + \mathbf{G}_i \mathbf{A}_i \hat{\mathbf{f}}_{m,i} \right) \right) \quad (17)$$

where  $\rho_i = (i + \epsilon)^{-r}$  is a mixing coefficient that controls the update rate,  $\pi_i = \frac{P}{|\mathbf{P}_i|}$  weights each update according to sample size,  $\epsilon$  is a delay hyperparameter and  $r$  is a forgetting rate (Hoffman et al. 2013). For the inverse scale,  $q(s)$  is updated using Equation 14, then its expected value is computed as  $\mathbb{E}[s] = \frac{\alpha}{\beta}$ .

The complete SVI algorithm is summarized in Algorithm 1. The use of an inner loop to learn  $\mathbf{G}_i$  avoids the need to store the complete matrix,  $\mathbf{G}$ . The inferred

**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  $G_i$  not converged do
    - 4 Compute  $G_i$  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  $\mathbb{E}[\ln s]$ ;
- 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.

distribution over the inducing points can be used for predicting the values of test items,  $f(\mathbf{x}^*)$ :

$$\mathbf{f}^* = \mathbf{K}_{*m} \mathbf{K}_{mm}^{-1} \hat{\mathbf{f}}_m, \quad \mathbf{C}^* = \mathbf{K}_{**} + \mathbf{K}_{*m} \mathbf{K}_{mm}^{-1} (\mathbf{S} - \mathbf{K}_{mm} / \mathbb{E}[s]) \mathbf{K}_{*m}^T \mathbf{K}_{mm}^{-1}, \quad (18)$$

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 items.

#### 4.3 SVI for CrowdGPPL

We now extend the SVI method to the crowd preference learning model proposed in Section 3.2. To begin with, we extend the variational posterior in Equation 11 to approximate the crowd model defined in Equation 6:

$$p(\mathbf{V}, \mathbf{V}_m, \mathbf{W}, \mathbf{W}_m, \mathbf{t}, \mathbf{t}_m, s_1, \dots, s_C, s_t | \mathbf{y}, \mathbf{x}, \mathbf{x}_m, \mathbf{u}, \mathbf{u}_m, k, \alpha_0, \beta_0) \quad (19)$$

$$\approx q(\mathbf{t})q(\mathbf{t}_m) \prod_{c=1}^C q(\mathbf{v}_c)q(\mathbf{w}_c)q(\mathbf{v}_{c,m})q(\mathbf{w}_{c,m})q(s_c)q(s_t) = q(\mathbf{F})q(s_t) \prod_{c=1}^C q(s_c),$$

where  $\mathbf{u}_m$  are the feature vectors of inducing users. By updating each of the  $q$  factors in turn, the SVI procedure optimizes the lower bound on the marginal likelihood (see also Equation 32 in the Appendix):

$$\begin{aligned} \mathcal{L}_{cr} = & \mathbb{E}_{q(\mathbf{F})} [\ln p(\mathbf{y} | \mathbf{F})] + \mathbb{E}_{q(\mathbf{t}_m), q(s_t)} [\ln p(\mathbf{t}_m, s_t | \mathbf{K}_{mm}, \alpha_0, \beta_0) - \ln q(\mathbf{t}_m)] \\ & \sum_{c=1}^C \left\{ \mathbb{E}_{q(\mathbf{v}_{m,c}), q(s_c)} [\ln p(\mathbf{v}_{m,c}, s_c | \mathbf{K}_{mm}, \alpha_0, \beta_0) - \ln q(\mathbf{v}_{m,c})] \right. \\ & \left. + \mathbb{E}_{q(\mathbf{w}_{m,c})} [\ln p(\mathbf{w}_{m,c} | \mathbf{L}_{mm}) - \ln q(\mathbf{w}_{m,c})] \right\}. \end{aligned} \quad (20)$$

The algorithm follows the same pattern as Algorithm 1, computing means and covariances for  $\mathbf{V}_m$ ,  $\mathbf{W}_m$  and  $\mathbf{T}_m$  instead of  $\mathbf{f}_m$ . The variational factor for the

inducing item factors is:

$$\begin{aligned}\ln q(\mathbf{v}_{m,c}) &= \ln \mathcal{N}(\mathbf{y}; \tilde{\Phi}(\mathbf{z}), Q) + \ln \mathcal{N}(\mathbf{v}_{m,c}; \mathbf{0}, \mathbf{K}_{mm}/\mathbb{E}[s_c]) + \text{const} \\ &= \ln \mathcal{N}(\mathbf{v}_{m,c}; \hat{\mathbf{v}}_{m,c}, \mathbf{S}_{v,c}),\end{aligned}\quad (21)$$

where the posterior mean  $\hat{\mathbf{v}}_{m,c}$  and covariance  $\mathbf{S}_{v,c}$  are computed using update equations of the same form as those of the single user GPPL in Equations 16 and 17. For reasons of space, we provide the complete equations for  $\hat{\mathbf{v}}_{m,c}$  and  $\mathbf{S}_{v,c}$  in Appendix B, Equations 33 and 34. The variational component for the inducing points of the common item mean follows a similar pattern:

$$\begin{aligned}\ln q(\mathbf{t}_m) &= \ln \mathcal{N}(\mathbf{y}; \tilde{\Phi}(\mathbf{z}), Q) + \ln \mathcal{N}(\mathbf{t}_m; \mathbf{0}, \mathbf{K}_{mm}/\mathbb{E}[s_t]) + \text{const} \\ &= \ln \mathcal{N}(\mathbf{t}_m; \hat{\mathbf{t}}_m, \mathbf{S}_t),\end{aligned}\quad (22)$$

where again, the posterior mean  $\hat{\mathbf{t}}$  and covariance  $\mathbf{S}_t$  use similar updates to Equation Equations 16 and 17, and their full definitions are given in Appendix B, Equations 36 and 37. Finally, the variational distribution for the inducing users factors is:

$$\begin{aligned}\ln q(\mathbf{w}_{m,c}) &= \ln \mathcal{N}(\mathbf{y}; \tilde{\Phi}(\mathbf{z}), Q) + \ln \mathcal{N}(\mathbf{w}_{m,c}; \mathbf{0}, \mathbf{L}_{mm}) + \text{const} \\ &= \ln \mathcal{N}(\mathbf{w}_{m,c}; \hat{\mathbf{w}}_{m,c}, \mathbf{\Sigma}_c),\end{aligned}\quad (23)$$

where  $\hat{\mathbf{w}}_c$  and  $\mathbf{\Sigma}_c$  also follow the pattern of Equations 16 and 17, and their full definitions are given in Appendix B, Equations 36 and 37. The expectations for the inverse scales,  $s_1, \dots, s_c$  and  $s_t$ , can be computed using Equation 14 by substituting in the corresponding terms for each  $\mathbf{v}_c$  or  $\mathbf{t}$  instead of  $\mathbf{f}$ .

Predictions for crowdGPPL can be made by computing the posterior mean utilities,  $\mathbf{F}^*$ , and the covariance  $\mathbf{A}_u^*$  for each user,  $u$ , in the test set:

$$\mathbf{F}^* = \hat{\mathbf{t}}^* + \sum_{c=1}^C \hat{\mathbf{v}}_c^{*T} \hat{\mathbf{w}}_c^*, \quad \mathbf{A}_u^* = \mathbf{C}_t^* + \sum_{c=1}^C \omega_{c,u}^* \mathbf{C}_{v,c}^* + \hat{w}_{c,u}^2 \mathbf{C}_{v,c}^* + \omega_{c,u}^* \hat{\mathbf{v}}_c \hat{\mathbf{v}}_c^T, \quad (24)$$

where  $\hat{\mathbf{t}}^*$ ,  $\hat{\mathbf{v}}_c^*$  and  $\hat{\mathbf{w}}_c^*$  are posterior means of the predictions,  $\mathbf{C}_t^*$  and  $\mathbf{C}_{v,c}^*$  are posterior item covariances of the predictions, and  $\omega_{c,u}^*$  is the posterior variance for the user components for the predicted users. These terms are defined in Appendix C, Equations 40 to 42. The mean  $\mathbf{F}^*$  and covariances  $\mathbf{A}_u^*$  can be inserted into Equations 3 and 2 to predict pairwise labels. In practice, the full covariance terms are needed only for Equations 3, so need only be computed between items for which we wish to predict pairwise labels.

## 5 Experiments

We test the key aspects of our proposed methods using the datasets summarized in Table 1: recovering an underlying consensus from noisy pairwise labels in a subjective task; modeling personal preferences from pairwise labels; and the scalability of our proposed Bayesian preference learning methods, GPPL and crowdGPPL

Page 12, Section 5. Can you provide a web reference to data sources (not just papers that use them)?

Dataset	#folds/ samples	#users	#items	train #pairs	#pairs	test #scores	#features items	users
Simulation a	25	25	100	900	0	100	2	2
Simulation b	25	25	100	900	0	100	2	2
Simulation c	25	25	100	36-2304	0	100	2	2
Sushi A-small	25	100	10	500	2500	1000	18	123
Sushi A	25	100	10	2000	2500	1000	18	123
Sushi B	25	5000	100	50000	5000	500000	18	123
UKPConvArg- CrowdSample	32	1442	1052	16398	529	33	32310	0

**Table 1** Summary of datasets showing mean counts per subsample or per fold. For the simulation datasets, we generate the subsamples of data independently, for Sushi we select subsamples independently from the dataset. UKPConvArgCrowdSample is divided into folds, where the test data in each fold corresponds to a single topic and stance. The numbers of features are given after categorical labels have been converted to one-hot encoding, counting each category as a separate feature.

using SVI. In Section 5.1, we use simulated data to test the robustness of our method to noise, small training sets and unknown numbers of latent components. Section 5.3 evaluates the ability of our approach to recover personal and consensus utilities on an NLP task with high-dimensional feature vectors. GPPL previously established the best performance on this dataset in Simpson and Gurevych (2018). Using the same dataset, we then analyze the scalability of our SVI approach. In Section 5.2, we compare our method against previous approaches for predicting the personal preferences of thousands of users on the *Sushi* datasets (Kamishima 2003). Finally, Section 5.4 evaluates whether the Bayesian approach was able to ignore redundant components when the number of components,  $C$ , is larger than necessary.

### 5.1 Simulated Noisy Data

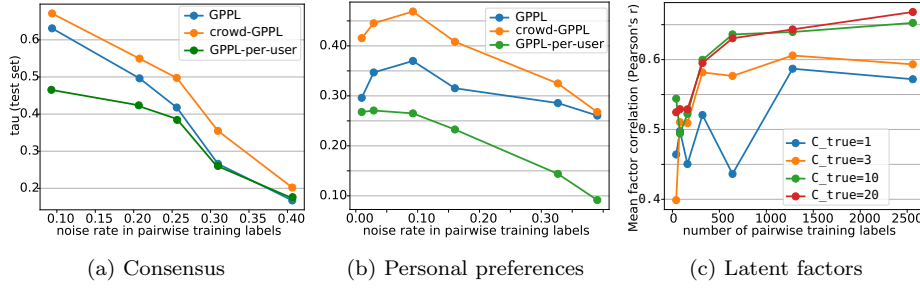
First, we test how well crowdGPPL is able to recover an underlying consensus function from pairwise labels with varying amounts of noise. We compare crowdGPPL against *GPPL* trained on all users’ preference labels to learn a single function, and separate GPPL instances per user with no collaborative learning (*GPPL-per-user*). The consensus for GPPL-per-user is the mean of all users’ predicted utilities. For crowdGPPL, we set  $C = 20$  and for all models, we set  $\alpha_0 = 1$ ,  $\beta_0 = 100$ , and use Matérn 3/2 kernels with length-scales chosen by a median heuristic:

$$l_{d,MH} = \text{median}(\{|x_{i,d} - x_{j,d}| \mid \forall i = 1, \dots, N, \forall j = 1, \dots, N\})D. \quad (25)$$

This is a computationally frugal heuristic for choosing the length-scales, which normalizes features and prevents the average covariance between items or users from shrinking as the number of features grows. The SVI hyperparameters were set to  $r = 0.9$ ,  $|P_i| = 1000$  and  $\epsilon = 1$ .

For the first test, we generate data by creating a  $20 \times 20$  grid of points and choosing 400 pairs of these points at random, which are split into 50% training and test sets. For each point, we generate pairwise labels by drawing from crowdGPPL

is D actually defined? I think the comment below stems from confusion that this is not just a median function. But anyway, why not add the refs for using a median? Page 13,



**Fig. 1** Simulations: Rank correlation between true and inferred preference values for different inference tasks. (a) & (b) varying level of noise in pairwise training labels, (c) varying number of pairwise training labels.

with 20 users, 5 latent components, and  $s = 0.001$ . We vary the precision of the consensus function,  $\sigma$ , to control the noise in the consensus function. The complete experiment was repeated 25 times, including generating new data for each value of  $\sigma$ . Figure 1a shows that the crowdGPPL model is better able to recover the latent consensus function than the other methods, even when noise levels are high. GPPL's predictions may be worsened by biased users whose preferences deviate consistently from the consensus. GPPL-per-user relies on separate instances of GPPL, so does not benefit from sharing information between users when training.

The second simulation modifies the previous setup by fixing  $\sigma = 10$  and varying  $s$  to evaluate the methods' ability to recover the personal preferences of simulated users. Results in Figure 1b show that crowdGPPL is able to make better predictions when noise is below 0.3 but its benefit disappears when the noise level increases further.

We hypothesize that learning a model with a larger number of latent components requires more training data. In the third simulation, we generate data using the same setup as before, but fix  $s = 0.2$  and  $\sigma = 1$  and vary the number of pairwise training labels and the number of true components through  $C_{true} \in \{1, 3, 10, 20\}$ . To evaluate the correlation between inferred and true user components, we compute Pearson correlations between each true component and each inferred component, then repeatedly select pairs of components with the highest correlations until every true component is matched to an inferred component. In Figure 1c we plot the mean of the correlations between matched pairs of components. For all values of  $C_{true}$ , increasing the number of training labels beyond 1000 leads to only minor increases in correlation at best. Performance is highest when  $C_{true} = 20$ , possibly because the predictive model has  $C = 25$  and hence is a closer match to the generating model. However, for all values of  $C_{true}$ , performance with  $> 500$  labels remains above 0.5, showing the model is reasonably robust to mismatches between  $C$  and  $C_{true}$ .

## 5.2 Sushi Preferences

We use sushi preference data (shown in Table 1), to compare GPPL and crowdGPPL with different sets of features against the previous benchmark method, *collabora-*

Page 14,  
Figure 1:  
include  
error  
bars, es-  
pecially  
in sub-  
figure  
(c).

*tive GP (collabGP)* (Houlsby et al. 2012). Collaborative GP uses a simpler model than crowdGPPL, with no consensus function, and no learning of the output scales of the utility functions, which makes it more important to select the correct number of factors in advance. It also uses an inference scheme based on expectation propagation (EP) and variational Bayes, which has a couple of important practical drawbacks: firstly, it does not infer the utilities, but learns pairwise labels directly, so its predictions cannot directly be used to rank or rate large sets of items; secondly, the number of parameters in the model is proportional to the dataset size, as it represents every combination of (user, item-pair) with a separate set of parameters, rather than working with inducing points, which may prevent the use of the model with large datasets due to memory constraints. However, as the EP approach of collabGP has been consistently shown to perform well for non-Gaussian likelihoods(?), we consider its performance as a target for our SVI approach for crowdGPPL.

The sushi datasets contain, for each user, a gold standard preference ranking of 10 types of sushi, from which we generate gold-standard pairwise labels. To test performance with very few training pairs, we obtain *Sushi-A-small* by selecting 100 users at random from the complete *Sushi-A* dataset, then selecting 5 pairs for training and 25 for testing per user. For *Sushi-A*, we select a subset of 100 users at random from the complete dataset, then split the data into training and test sets by randomly selecting 20 pairs for each user for training and 25 for testing. For *Sushi-B*, we use all 5000 workers, and subsample 10 training pairs and 1 test pair per user. Beside GPPL, crowdGPPL and GPPL-per-user, we introduce four further baselines: *crowdGPPL*\u, which ignores the user features; *crowdGPPL*\u\x, which ignores both user and item features; *crowdGPPL*\u\t, which excludes the consensus function  $t$  from the model as well as the user features; and *crowdGPPL*\induc, which uses all features but does not use inducing points for a sparse approximation. For \u methods, the user covariance matrix,  $\mathbf{K}_w$ , in the crowdGPPL model is replaced by the identity matrix, and for crowdGPPL\u\x, the item covariance matrices,  $\mathbf{K}_v$  and  $\mathbf{K}_t$  are also replaced by the identity matrix. We set hyperparameters  $C = 20$  without optimization,  $\alpha_0 = 1, \beta_0 = 100$  using a grid search over values  $10^{\{-1, \dots, 3\}}$  on withheld user data from the *Sushi-A* dataset,  $\epsilon = 5$ ,  $|P_i| = 200$  for *Sushi-A* and  $|P_i| = 2000$  for *Sushi-B*. All other hyperparameters are the same as for Section 5.1. The complete process of subsampling, training and testing, was repeated 25 times for each dataset.

The results in Table 2 illustrate the performance benefit of crowd models over single-user GPPL, and the runtimes show the speedup of the sparse approximation of crowdGPPL over crowdGPPL\induc. GPPL is substantially quicker to train than crowdGPPL, and GPPL-per-user does not scale well to larger numbers of users. When using inducing points, the user features decrease the performance of crowdGPPL: crowdGPPL\induc and crowdGPPL\u both outperform the full crowdGPPL. To use inducing points, there must be a strong relationship between neighbouring points, which may not to be the case given the features in this dataset, since they describe only very general characteristics, such as the user's region and age group. Comparing crowdGPPL\u with crowdGPPL\u\t, including the consensus function appears to improve performance by a modest amount. The strong performance of GPPL-per-user suggests that even ten preferences per person were enough to learn a reasonable model for *Sushi-B*.

some  
comment  
on the  
strength  
of  
GPPL-  
per-user?

does  
houlsby  
learn the  
output  
scales of  
the fac-  
tors? If  
not, then  
they  
are not

Method	Sushi-A-small			Sushi-A			Sushi-B		
	Acc	CEE	$\tau$	Acc	CEE	$\tau$	Acc	CEE	$\tau$
crowdGPPL	.68	.03	.39	.81	.39	.66	.73	.57	.47
crowdGPPL \inducing	.70	<b>.57</b>	.46	.84	.33	.79	-	-	-
crowdGPPL \u	.70	.58	.46	<b>.85</b>	<b>.31</b>	<b>.81</b>	.78	<b>.57</b>	.49
crowdGPPL \u, $C = 50$				.83	.34	.78			
crowdGPPL \u\x	<b>.71</b>	<b>.57</b>	<b>.49</b>	<b>.85</b>	.33	.80	.77	.58	.48
crowdGPPL \u,t	.68	.60	.43	.84	.33	.80	.76	.61	.54
GPPL	.65	.62	.31	.65	.62	.31	.65	.62	.31
GPPL-per-user	.67	.64	.42	.83	.40	.79	.75	.60	<b>.60</b>
collabGP	.69	.58	n/a	.83	.35	n/a	.76	.49	n/a

**Table 2** Predicting personal preferences: performance on *Sushi-A* dataset and *Sushi-B* datasets, means over 25 repeats. CrowdGPPL uses  $C = 20$  unless otherwise specified. The standard deviations of all metrics shown are  $\leq 0.02$ , for CEE,  $\leq 0.08$ , for  $\tau$ ,  $\leq 0.03$ . For Sushi-B, crowdGPPL, GPPL-per-user and collabGP had runtimes of approximately 30 minutes on a 12 core, 2.6GHz CPU server, while GPPL required only one minute.

Method	Consensus			Personal: all workers			>40 training pairs		
	Acc	CEE	$\tau$	Acc	CEE	$\tau$	Acc	CEE	$\tau$
GPPL	.77	.50	.40	.71	.56	.32	.72	.55	.26
crowdGPPL	.78	.50	.53	.73	.61	.33	.74	.59	.29
crowdBT-GP	.77	.52	.50	.71	.60				

**Table 3** UKPConvArgCrowdSample, using ling+GloVe features: predicting consensus values, personal preferences, and personal preferences for workers with  $\geq 40$  pairs in the training set. Classification accuracy (*Acc*) and cross entropy error (or log-loss, *CEE*) evaluate pairwise predictions, Kendall’s  $\tau$  evaluates the predicted utilities. Runtimes on a 12 core, 2.6GHz CPU server were approximately 3 minutes for GPPL and crowdBT-GP, 60 minutes for crowdGPPL.

Like crowdGPPL, collaborative GP (collabGP) considers user features as inputs but requires many more model parameters. This allows it to outperform crowdGPPL on accuracy and CEE but does not provide a ranking function for computing Kendall’s  $\tau$ . The inference method of collabGP did not present practical problems on the Sushi dataset on computer with 16GB RAM, and in fact produces relatively fast run times due to its more efficient implementation using a C library in place of Python. However, the results show that crowdGPPL is able to obtain competitive performance despite the approximations required for our proposed SVI method. In the next experiment, we test the approach on a larger dataset with many more pairs, users and items, for which the memory requirements of collabGP become problematic.

### 5.3 Argument Convincingness

We evaluate consensus learning, personal preference learning and scalability on an NLP task, namely identifying convincing arguments. The dataset, *UKPConvArgCrowdSample*, was subsampled by Simpson and Gurevych (2018) from the crowdsourced data provided by Habernal and Gurevych (2016), and contains arguments written by users of online debating forums, along with crowdsourced judgments of pairs of arguments indicating the most convincing argument. Each argument is represented by 32,310 numerical features and the dataset is divided

state why we couldn't use collabGP for convincingness, or run it and add it to the results. Consider why Khan does not compare against collabGP? And check that the memory requirements thing is the core reason



into 32 folds (16 topics, each of which has two opposing stances). For each fold, we train on 31 folds and test on the remaining fold. We extend the task described in Simpson and Gurevych (2018) to predict not just the consensus, but also the personal preferences of individual crowd workers. GPPL was previously shown to outperform SVM, Bi-LSTM and Gaussian process classifier methods at consensus prediction for *UKPConvArgCrowdSample* (Simpson and Gurevych 2018). We hypothesize that a worker’s view of convincingness depends on their prior beliefs and understanding of the subject discussed, and that crowdGPPL may therefore predict unseen pairwise labels or rankings for individual workers or the consensus more accurately than GPPL, by accounting for the biases of individual workers.

Table 3 shows performance for GPPL and crowdGPPL. The hyperparameters were kept the same as in Section 5.1 except for: GPPL uses  $\alpha_0 = 2$ ,  $\beta_0 = 200$  and crowdGPPL uses  $\alpha_0 = 2$ ,  $\beta_0 = 2000$ , set by comparing training set accuracy against  $\alpha_0 = 2, \beta_0 = 20000$  and  $\alpha_0 = 2, \beta_0 = 2$ ;  $C = 50$  and  $\epsilon = 2$ , which were not optimized. CrowdGPPL outperforms GPPL at predicting the consensus pairwise labels shown by classification accuracy (*acc*) and cross entropy error (*CEE*), and the consensus ranking (significant with  $p < 0.05$ , Wilcoxon signed-rank test), shown by Kendall’s  $\tau$  rank correlation. For the personal preference predictions, crowdGPPL also outperforms GPPL at ranking pairwise label accuracy, suggesting that there is a benefit to modeling individual workers when predicting the consensus. The benefits of crowdGPPL on this task may be restricted because the pairwise comparisons in the test folds are noisy and contain numerous contradictions (Habernal and Gurevych 2016). For workers with more training data, the more complex crowdGPPL model is able to improve further over GPPL. Since the CEE scores of crowdGPPL are lower than those of GPPL, while accuracy is higher, crowdGPPL may be slightly under-confident.

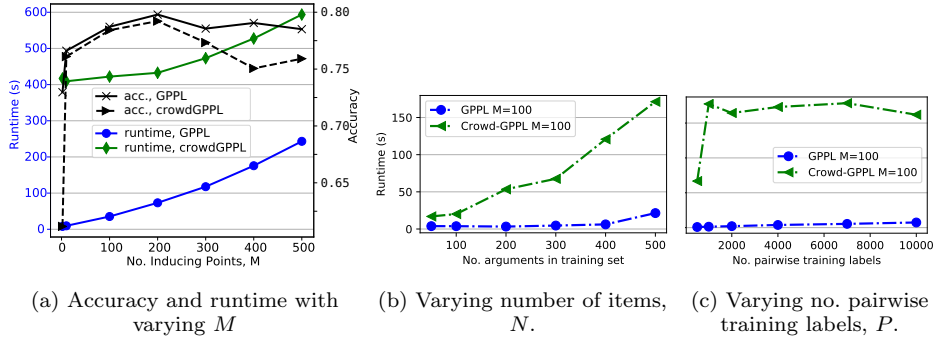
We examine the scalability of our SVI implementation by evaluating GPPL and crowdGPPL with different numbers of inducing points,  $M$ . Here, we fix  $C = 5$  and keep other model hyperparameters and experimental setup the same as for Table 3. Figure 2a shows the trade-off between runtime and accuracy as an effect of choosing  $M$ . Accuracy peaks using  $M = 200$  while the runtime continues to increase rapidly in a polynomial fashion. Since there are 33,210 features, the runtime includes large overheads due to the computation of the covariance matrices, which is linear in the number of features.

Figures 2b and 2c show runtimes as a function of the number of items in the training set,  $N_{tr}$ , and the number of pairwise training labels,  $P$ , respectively (all other settings remain as in Figure 2a). The runtime increase with  $N_{tr}$  for GPPL is almost imperceptible, while for crowdGPPL it increases almost linearly as more computations over the set of items are required than for GPPL, and the method takes more iterations to converge with more items. However, many of the additional computations can be parallelized in future implementations. Increasing the number of pairwise labels,  $P$ , above 1000, however, does not visibly affect the runtimes.

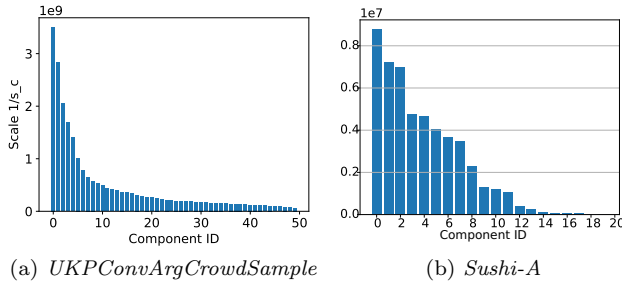
#### 5.4 Posterior Variance of Item Components

We investigate how many latent components were actively used by crowdGPPL on the *UKPConvArgCrowdSample* and *Sushi-A* datasets using the median heuristic.

do we  
say that  
we use  
m=500  
and  
P=200?



**Fig. 2** Runtimes for training+prediction on UKPConvArgCrowdSample. 300 GloVe features. Means over 15 runs. CrowdGPPL with 5 components.



**Fig. 3** Distribution of latent factor variances,  $1/s_c$ , for crowdGPPL on UKPConvArgCrowdSample, *Sushi-A* and *Sushi-B*, averaged over all runs.

Figure 3 plots the posterior expectations of the inferred scales,  $1/s_c$ , for the latent item components. The plots show that many factors have a very small variance and therefore do not contribute strongly to many of the model’s predictions. This indicates that our Bayesian approach, in which the priors of the latent factors have mean zero, has inferred a simpler model than the number of latent components would permit.

## 6 Conclusions

We proposed a novel Bayesian preference learning approach for jointly modeling the individual utilities of members of a crowd and forming a consensus from unreliable annotations, such as those provided by a crowd. Unlike previous methods, our model learns the latent utilities of items from pairwise comparisons using a combination of Gaussian processes and Bayesian matrix factorization to capture divergences in opinion while inferring a consensus. To enable inference at the scale of real-world datasets in fields such as NLP, we derived a stochastic variational inference scheme. Our empirical results confirm that our approach scales well with the amount of training data, and that jointly modeling the consensus and personal preferences of users can improve the predictions of both. When predicting a con-

sensus from crowdsourced data, our model, CrowdGPPL, improves on the results of Simpson and Gurevych (2018), who used a GPPL method that did not account for personal preferences. On a benchmark recommendation dataset, our approach also produces competitive results and showed that modeling the consensus function can boost predictions of individual preferences. Future work will evaluate the benefit of learning inducing point locations from data and optimizing length-scale hyperparameters by maximizing  $\mathcal{L}$ .

## 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: Twenty-Third International Joint Conference on Artificial Intelligence, pp 1213–1219
- 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
- Banerji M, Lahav O, Lintott CJ, Abdalla FB, Schawinski K, Bamford SP, Andreescu D, Murray P, Raddick MJ, Slosar A, et al. (2010) Galaxy Zoo: reproducing galaxy morphologies via machine learning. *Monthly Notices of the Royal Astronomical Society* 406(1):342–353
- Bradley RA, Terry ME (1952) Rank analysis of incomplete block designs: I. the method of paired comparisons. *Biometrika* 39(3/4):324–345
- Chen G, Zhu F, Heng PA (2018) Large-scale Bayesian probabilistic matrix factorization with memo-free distributed variational inference. *ACM Transactions on Knowledge Discovery from Data* 12(3):31:1–31:24
- 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
- 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
- Fürnkranz J, Hüllermeier E (2010) Preference learning and ranking by pairwise comparison. In: *Preference learning*, Springer, pp 65–82
- Gaunt A, Borsa D, Bachrach Y (2016) Training deep neural nets to aggregate crowdsourced responses. In: Proceedings of the Thirty-Second Conference on Uncertainty in Artificial Intelligence. AUAI Press, p 242251
- Guo S, Sanner S, Bonilla EV (2010) Gaussian process preference elicitation. In: *Advances in neural information processing systems*, pp 262–270
- Habernal I, Gurevych I (2016) Which argument is more convincing? Analyzing and predicting convincingness of Web arguments using bidirectional LSTM. In: Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), Association for Computational Linguistics, pp 1589–1599
- 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
- Joachims T (2002) Optimizing search engines using clickthrough data. In: *Proceedings of the eighth ACM SIGKDD international conference on Knowledge discovery and data mining*, ACM, pp 133–142
- Kamishima T (2003) Nantonac collaborative filtering: recommendation based on order responses. In: *Proceedings of the ninth ACM SIGKDD international conference on Knowledge discovery and data mining*, ACM, pp 583–588
- Kendall MG (1948) Rank correlation methods. Griffin
- Khan ME, Ko YJ, Seeger M (2014) Scalable collaborative Bayesian preference learning. In: *Proceedings of the 17th International Conference on Artificial Intelligence and Statistics*, vol 33, 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
- Kingsley DC, Brown TC (2010) Preference uncertainty, preference refinement and paired comparison experiments. *Land Economics* 86(3):530–544
- Kiritchenko S, Mohammad SM (2016) Capturing reliable fine-grained sentiment associations by crowdsourcing and best–worst scaling. In: *Proceedings of NAACL-HLT*, pp 811–817
- Koren Y, Bell R, Volinsky C (2009) Matrix factorization techniques for recommender systems. *Computer* (8):30–37
- Li J, Baba Y, Kashima H (2018) Simultaneous clustering and ranking from pairwise comparisons. In: *Proceedings of the 27th International Joint Conference on Artificial Intelligence*, AAAI Press, pp 1554–1560
- Luce RD (1959) On the possible psychophysical laws. *Psychological review* 66(2):81
- Lukin S, Anand P, Walker M, Whittaker S (2017) Argument strength is in the eye of the beholder: Audience effects in persuasion. In: *Proceedings of the 15th Conference of the European Chapter of the Association for Computational Linguistics*, pp 742–753
- 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
- Ovadia S (2004) Ratings and rankings: Reconsidering the structure of values and their measurement. *International Journal of Social Research Methodology* 7(5):403–414
- Plackett RL (1975) The analysis of permutations. *Applied Statistics* pp 193–202
- 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
- Resnick P, Varian HR (1997) Recommender systems. *Communications of the ACM* 40(3):56–58
- 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
- Salimans T, Paquet U, Graepel T (2012) Collaborative learning of preference rankings. In: *Proceedings of the sixth ACM conference on Recommender systems*, ACM, pp 261–264
- 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
- Simpson E, Gurevych I (2018) Finding convincing arguments using scalable bayesian preference learning. *Transactions of the Association for Computational Linguistics* 6:357–371

- Simpson E, Reece S, Roberts SJ (2017) Bayesian heatmaps: probabilistic classification with multiple unreliable information sources. In: Joint European Conference on Machine Learning and Knowledge Discovery in Databases, Springer, pp 109–125
- Snelson E, Ghahramani Z (2006) Sparse Gaussian processes using pseudo-inputs. In: Advances in neural information processing systems, pp 1257–1264
- Snow R, O’Connor B, Jurafsky D, Ng AY (2008) Cheap and fast—but is it good?: evaluating non-expert annotations for natural language tasks. In: Proceedings of the conference on empirical methods in natural language processing, Association for Computational Linguistics, pp 254–263
- 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
- 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
- Yang YH, Chen HH (2011) Ranking-based emotion recognition for music organization and retrieval. *IEEE Transactions on Audio, Speech, and Language Processing* 19(4):762–774
- Yannakakis GN, Hallam J (2011) Ranking vs. preference: a comparative study of self-reporting. In: International Conference on Affective Computing and Intelligent Interaction, Springer, pp 437–446
- 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

## A Deriving the Variational Lower Bound for GPPL

Due to the non-Gaussian likelihood, Equation 2, the posterior distribution over  $\mathbf{f}$  contains intractable integrals:

$$p(\mathbf{f}|\mathbf{y}, k_\theta, \alpha_0, \alpha_0) = \frac{\int \prod_{p=1}^P \Phi(z_p) \mathcal{N}(\mathbf{f}; \mathbf{0}, \mathbf{K}_\theta/s) \mathcal{G}(s; \alpha_0, \beta_0) ds}{\int \int \prod_{p=1}^P \Phi(z_p) \mathcal{N}(\mathbf{f}'; \mathbf{0}, \mathbf{K}_\theta/s) \mathcal{G}(s; \alpha_0, \beta_0) ds d\mathbf{f}'} \quad (26)$$

We can derive a variational lower bound as follows, beginning with an approximation that does not use inducing points:

$$\begin{aligned} \mathcal{L}_1 = \sum_{p=1}^P \mathbb{E}_{q(\mathbf{f})} \left[ \ln p(y(a_p, b_p) | f(\mathbf{x}_{a_p}), f(\mathbf{x}_{b_p})) \right] &+ \mathbb{E}_{q(\mathbf{f}), q(s)} \left[ \ln \frac{p(\mathbf{f} | \mathbf{0}, \mathbf{K}/s)}{q(\mathbf{f})} \right] \\ &+ \mathbb{E}_{q(s)} \left[ \ln \frac{p(s | \alpha_0, \beta_0)}{q(s)} \right] \end{aligned} \quad (27)$$

Substituting the forms of the distributions with their variational parameters, we get:

$$\begin{aligned} \mathcal{L}_1 = \mathbb{E}_{q(\mathbf{f})} \left[ \sum_{p=1}^P y(a_p, b_p) \ln \Phi(z_p) + y(b_p, a_p) (1 - \ln \Phi(z_p)) \right] \\ + \ln \mathcal{N}(\hat{\mathbf{f}}; \boldsymbol{\mu}, \mathbf{K}/\mathbb{E}[s]) - \ln \mathcal{N}(\hat{\mathbf{f}}; \hat{\mathbf{f}}, \mathbf{C}) + \mathbb{E}_{q(s)} [\ln \mathcal{G}(s; \alpha_0, \beta_0) - \ln \mathcal{G}(s; \alpha, \beta)] \end{aligned} \quad (28)$$

We now replace the likelihood with a Gaussian approximation:

$$\begin{aligned}
\mathcal{L}_1 &\approx \mathcal{L}_2 = \mathbb{E}_{q(\mathbf{f})} [\mathcal{N}(\mathbf{y}|\Phi(\mathbf{z}), \mathbf{Q})] + \ln \mathcal{N}(\hat{\mathbf{f}}; \boldsymbol{\mu}, \mathbf{K}/\mathbb{E}[s]) - \ln \mathcal{N}(\hat{\mathbf{f}}; \hat{\mathbf{f}}, \mathbf{C}) \\
&\quad + \mathbb{E}_{q(s)} [\ln \mathcal{G}(s; \alpha_0, \beta_0) - \ln \mathcal{G}(s; \alpha, \beta)] \\
&= -\frac{1}{2} \left\{ L \ln 2\pi + \ln |\mathbf{Q}| - \ln |\mathbf{C}| + \ln |\mathbf{K}| - \mathbb{E}[\ln s] + (\hat{\mathbf{f}} - \boldsymbol{\mu})\mathbb{E}[s]\mathbf{K}^{-1}(\hat{\mathbf{f}} - \boldsymbol{\mu}) \right. \\
&\quad \left. + \mathbb{E}_{q(\mathbf{f})} \left[ (\mathbf{y} - \Phi(\mathbf{z}))^T \mathbf{Q}^{-1} (\mathbf{y} - \Phi(\mathbf{z})) \right] \right\} - \Gamma(\alpha_0) + \alpha_0(\ln \beta_0) + (\alpha_0 - \alpha)\mathbb{E}[\ln s] \\
&\quad + \Gamma(\alpha) + (\beta - \beta_0)\mathbb{E}[s] - \alpha \ln \beta,
\end{aligned} \tag{29}$$

where  $\mathbb{E}[s] = \frac{\alpha}{\beta}$ ,  $\mathbb{E}[\ln s] = \Psi(2\alpha) - \ln(2\beta)$ ,  $\Psi$  is the digamma function and  $\Gamma(\cdot)$  is the gamma function. Finally, we use a Taylor-series linearisation to make the remaining expectation tractable:

$$\begin{aligned}
\mathcal{L}_2 &\approx \mathcal{L}_3 = -\frac{1}{2} \left\{ L \ln 2\pi + \ln |\mathbf{Q}| - \ln |\mathbf{C}| + \ln |\mathbf{K}/\mathbb{E}[s]| + (\hat{\mathbf{f}} - \boldsymbol{\mu})\mathbb{E}[s]\mathbf{K}^{-1}(\hat{\mathbf{f}} - \boldsymbol{\mu}) \right. \\
&\quad \left. + (\mathbf{y} - \Phi(\hat{\mathbf{z}}))^T \mathbf{Q}^{-1} (\mathbf{y} - \Phi(\hat{\mathbf{z}})) \right\} - \Gamma(\alpha_0) + \alpha_0(\ln \beta_0) + (\alpha_0 - \alpha)\mathbb{E}[\ln s] \\
&\quad + \Gamma(\alpha) + (\beta - \beta_0)\mathbb{E}[s] - \alpha \ln \beta.
\end{aligned} \tag{30}$$

Now, we can introduce the sparse approximation to obtain the bound in Equation 15:

$$\begin{aligned}
\mathcal{L} &\approx \mathcal{L}_3 = \mathbb{E}_{q(\mathbf{f})} [\ln p(\mathbf{y}|\mathbf{f})] + \mathbb{E}_{q(\mathbf{f}_m), q(s)} [\ln p(\mathbf{f}_m, s|\mathbf{K}, \alpha_0, \beta_0)] - \mathbb{E}_{q(\mathbf{f}_m)} [\ln q(\mathbf{f}_m)] - \mathbb{E}_{q(s)} [\ln q(s)] \\
&= \sum_{p=1}^P \mathbb{E}_{q(\mathbf{f})} [\ln p(y_p|f_{a_p}, f_{b_p})] - \frac{1}{2} \left\{ \ln |\mathbf{K}_{mm}| - \mathbb{E}[\ln s] - \ln |\mathbf{S}| - M \right. \\
&\quad \left. + \hat{\mathbf{f}}_m^T \mathbb{E}[s] \mathbf{K}_{mm}^{-1} \hat{\mathbf{f}}_m + \text{tr}(\mathbb{E}[s] \mathbf{K}_{mm}^{-1} \mathbf{S}) \right\} + \ln \Gamma(\alpha) - \ln \Gamma(\alpha_0) + \alpha_0(\ln \beta_0) \\
&\quad + (\alpha_0 - \alpha)\mathbb{E}[\ln s] + (\beta - \beta_0)\mathbb{E}[s] - \alpha \ln \beta,
\end{aligned} \tag{31}$$

where the terms relating to  $\mathbb{E}[p(\mathbf{f}|\mathbf{f}_m) - q(\mathbf{f})]$  cancel. For crowdGPPL, our approximate variational lower bound is:

$$\begin{aligned}
\mathcal{L}_{cr} &= \sum_{p=1}^P \ln p(y_p|\hat{\mathbf{v}}_{\cdot, a_p}^T \hat{\mathbf{w}}_{\cdot, j_p} + \hat{t}_{a_p}, \hat{\mathbf{v}}_{\cdot, b_p}^T \hat{\mathbf{w}}_{\cdot, j_p} + \hat{t}_{b_p}) - \frac{1}{2} \left\{ \sum_{c=1}^C \left\{ -M_n - M_u + \ln |\mathbf{K}_{mm}| \right. \right. \\
&\quad \left. \left. + \ln |\mathbf{L}_{mm}| - \ln |\mathbf{S}_{v,c}| - \mathbb{E}[\ln s_c] + \hat{\mathbf{v}}_{m,c}^T \mathbb{E}[s_c] \mathbf{K}_{mm}^{-1} \hat{\mathbf{v}}_{m,c} + \text{tr}(\mathbb{E}[s_c] \mathbf{K}_{mm}^{-1} \mathbf{S}_{v,c}) \right. \right. \\
&\quad \left. \left. - \ln |\mathbf{S}_c| + \hat{\mathbf{w}}_{m,c}^T \mathbf{L}_{mm}^{-1} \hat{\mathbf{w}}_{m,c} + \text{tr}(\mathbf{L}_{mm}^{-1} \mathbf{S}_c) \right\} - M_n + \ln |\mathbf{K}_{mm}| - \ln |\mathbf{S}_t| - \mathbb{E}[\ln s_t] \right. \\
&\quad \left. + \hat{\mathbf{t}}^T \mathbb{E}[s_t] \mathbf{K}_{mm}^{-1} \hat{\mathbf{t}} + \text{tr}(\mathbb{E}[s_t] \mathbf{K}_{mm}^{-1} \mathbf{S}_t) \right\} - (C+1)(\ln \Gamma(\alpha_0) + \alpha_0(\ln \beta_0)) \\
&\quad + \sum_{c=1}^C \left\{ \ln \Gamma(\alpha_c) + (\alpha_0 - \alpha_c)\mathbb{E}[\ln s_c] + (\beta_c - \beta_0)\mathbb{E}[s_c] - \alpha_c \ln \beta_c \right\} \\
&\quad + \ln \Gamma(\alpha_{s_t}) + (\alpha_0 - \alpha_{s_t})\mathbb{E}[\ln s_t] + (\beta_{s_t} - \beta_0)\mathbb{E}[s_c] - \alpha_{s_t} \ln \beta_{s_t},
\end{aligned} \tag{32}$$

## B Posterior Parameters for Variational Factors in CrowdGPPL

For the latent item components, the posterior precision estimate for  $\mathbf{S}_{v,c}^{-1}$  at iteration  $i$  is given by:

$$\mathbf{S}_{v,c,i}^{-1} = (1 - \rho_i) \mathbf{S}_{v,c,i-1}^{-1} + \rho_i \left( \mathbf{K}_{mm}^{-1} \mathbb{E}[s_c] + \pi_i \mathbf{A}_{v,i}^T \mathbf{G}_i^T \text{diag}(\hat{\mathbf{w}}_{c,u}^2 + \mathbf{S}_{c,u,u}) \mathbf{Q}_i^{-1} \mathbf{G}_i \mathbf{A}_{v,i} \right), \tag{33}$$

where  $\mathbf{A}_i = \mathbf{K}_{im}\mathbf{K}_{mm}^{-1}$ ,  $\hat{\mathbf{w}}_c$  and  $\Sigma_c$  are the variational mean and covariance of the  $c$ th latent user component (defined below in Equations 39 and 38), and  $\mathbf{u} = \{u_p \forall p \in P_i\}$  is the vector of user indexes in the sample of observations. We use  $\mathbf{S}_{v,c}^{-1}$  to compute the means for each row of  $\mathbf{V}_m$ :

$$\hat{\mathbf{v}}_{m,c,i} = \mathbf{S}_{v,c,i} \left( (1 - \rho_i) \mathbf{S}_{v,c,i-1}^{-1} \hat{\mathbf{v}}_{m,c,i-1} + \rho_i \pi_i \left( \mathbf{S}_{v,c,i} \mathbf{A}_i^T \mathbf{G}_i^T \text{diag}(\hat{\mathbf{w}}_c, \mathbf{u}) \mathbf{Q}_i^{-1} \right. \right. \\ \left. \left. \left( \mathbf{y}_i - \Phi(\mathbb{E}[\mathbf{z}_i]) + \sum_{j=1}^U \mathbf{H}_j^{(i)} (\hat{\mathbf{v}}_c^T \hat{\mathbf{w}}_{c,j}) \right) \right) \right), \quad (34)$$

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

$$H_{j,p,n}^{(i)} = \Phi(\mathbb{E}[\mathbf{z}_p]) (1 - \Phi(\mathbb{E}[\mathbf{z}_p])) (2y_p - 1) ([n = a_p] - [n = b_p]) [j = u_p]. \quad (35)$$

For the consensus, the precision and mean are updated according to the following:

$$\mathbf{S}_{t,i}^{-1} = (1 - \rho_i) \mathbf{S}_{t,i-1}^{-1} + \rho_i \mathbf{K}_{t,mm}^{-1} / \mathbb{E}[s_t] + \rho_i \pi_i \mathbf{A}_i^T \mathbf{G}_i^T \mathbf{Q}_i^{-1} \mathbf{G}_i \mathbf{A}_i \quad (36)$$

$$\hat{\mathbf{t}}_{m,i} = \mathbf{S}_{t,i} \left( (1 - \rho_i) \mathbf{S}_{t,i-1}^{-1} \hat{\mathbf{t}}_{m,i-1} + \rho_i \pi_i \mathbf{A}_i^T \mathbf{G}_i^T \mathbf{Q}_i^{-1} (\mathbf{y}_i - \Phi(\mathbb{E}[\mathbf{z}_i]) + \mathbf{G}_i \mathbf{A}_i \hat{\mathbf{t}}_i) \right). \quad (37)$$

For the latent user components, the SVI updates for the parameters are:

$$\Sigma_{c,i}^{-1} = (1 - \rho_i) \Sigma_{c,i-1}^{-1} + \rho_i \mathbf{L}_{mm}^{-1} + \rho_i \pi_i \mathbf{A}_{w,i}^T \left( \sum_{p \in P_i} \mathbf{H}_{\cdot,p}^{(i)T} \text{diag}(\hat{\mathbf{v}}_{c,a}^2 + \right. \\ \left. S_{c,a,a} + \hat{\mathbf{v}}_{c,b}^2 + S_{c,b,b} - 2\hat{\mathbf{v}}_{c,a} \hat{\mathbf{v}}_{c,b} - 2S_{c,a,b}) \mathbf{Q}_i^{-1} \sum_{p \in P_i} \mathbf{H}_{\cdot,p}^{(i)} \right) \mathbf{A}_{w,i} \quad (38)$$

$$\hat{\mathbf{w}}_{m,c,i} = \Sigma_{c,i} \left( (1 - \rho_i) \Sigma_{c,i-1} \hat{\mathbf{w}}_{m,c,i-1} + \rho_i \pi_i \mathbf{A}_{w,i}^T \sum_{p \in P_i} \mathbf{H}_{\cdot,p}^{(i)} (\text{diag}(\hat{\mathbf{v}}_{c,a}) \right. \\ \left. - \text{diag}(\hat{\mathbf{v}}_{c,b})) \mathbf{Q}_i^{-1} \left( \mathbf{y}_i - \Phi(\mathbb{E}[\mathbf{z}_i]) + \sum_{j=1}^U \mathbf{H}_j^{(i)} (\hat{\mathbf{v}}_c^T \hat{\mathbf{w}}_{c,j}) \right) \right), \quad (39)$$

where the subscripts  $\mathbf{a} = \{a_p \forall p \in P_i\}$  and  $\mathbf{b} = \{b_p \forall p \in P_i\}$  are lists of indices to the first and second items in the pairs, respectively, and  $\mathbf{A}_{w,i} = \mathbf{L}_{im} \mathbf{L}_{mm}^{-1}$ .

## C Predictions with CrowdGPPL

The means, item covariances and user variance required for predictions with crowdGPPL (Equation 24) are defined as follows:

$$\hat{\mathbf{t}}^* = \mathbf{K}_{*m} \mathbf{K}_{mm}^{-1} \hat{\mathbf{t}}_m, \quad \hat{\mathbf{v}}_c^* = \mathbf{K}_{*m} \mathbf{K}_{mm}^{-1} \hat{\mathbf{v}}_{m,c}, \quad \hat{\mathbf{w}}_c^* = \mathbf{L}_{*m} \mathbf{L}_{mm}^{-1} \hat{\mathbf{w}}_{m,c}, \quad (40)$$

$$\mathbf{C}_t^* = \frac{\mathbf{K}_{**}}{\mathbb{E}[s_t]} + \mathbf{A}_{*m} (\mathbf{S}_t - \mathbf{K}_{mm}) \mathbf{A}_{*m}^T, \quad \mathbf{C}_{v,c}^* = \frac{\mathbf{K}_{**}}{\mathbb{E}[s_c]} + \mathbf{A}_{*m} (\mathbf{S}_{v,c} - \mathbf{K}_{mm}) \mathbf{A}_{*m}^T \quad (41)$$

$$\omega_{c,u}^* = 1 + \mathbf{A}_{w,um} (\Sigma_{w,c} - \mathbf{L}_{mm}) \mathbf{A}_{w,um}^T \quad (42)$$

where  $\mathbf{A}_{*m} = \mathbf{K}_{*m} \mathbf{K}_{mm}^{-1}$ ,  $\mathbf{A}_{w,um} = \mathbf{L}_{um} \mathbf{L}_{mm}^{-1}$  and  $\mathbf{L}_{um}$  is the covariance between user  $u$  and the inducing users.