
Latent Variable Matrix Factorization in Heterogeneous Peer Prediction

Brian Hentschel Anna S. Hilgard Casey Meehan

Abstract

- This document describes the expected style, structure, and rough proportions for your final project write-up.
- While you are free to break from this structure, consider it a strong prior for our expectations of the final report.
- Length is a hard constraint. You are only allowed max **8 pages** in this format. While you can include supplementary material, it will not be factored into the grading process. It is your responsibility to convey the main contributions of the work in the length given.

Lorem ipsum dolor sit amet, consectetur adipiscing elit. Ut purus elit, vestibulum ut, placerat ac, adipiscing vitae, felis. Curabitur dictum gravida mauris. Nam arcu libero, nonummy eget, consectetur id, vulputate a, magna. Donec vehicula augue eu neque. Pellentesque habitant morbi tristique senectus et netus et malesuada fames ac turpis egestas. Mauris ut leo. Cras viverra metus rhoncus sem. Nulla et lectus vestibulum urna fringilla ultrices. Phasellus eu tellus sit amet tortor gravida placerat. Integer sapien est, iaculis in, pretium quis, viverra ac, nunc. Praesent eget sem vel leo ultrices bibendum. Aenean faucibus. Morbi dolor nulla, malesuada eu, pulvinar at, mollis ac, nulla. Curabitur auctor semper nulla. Donec varius orci eget risus. Duis nibh mi, congue eu, accumsan eleifend, sagittis quis, diam. Duis eget orci sit amet orci dignissim rutrum.

1. Introduction

In peer prediction settings with heterogeneous users, the goal is to make truthful reporting a low-regret strategy for any agent. To fully minimize expected regret while allowing for heterogeneity, this would require the computation and use of individualized user-user signal correlations for every pair of users. This is in general computationally infeasible, so methodologies typically focus instead on clus-

tering similar users while trying to minimally distort their individual signal distributions.

Additionally, because this method explicitly requires an existing signal distribution, there is a cold start problem: if we do not already have a significant amount of task data for a new user, we cannot know how properly to reward them in comparison to their peers. In our project, we focus on two specific goals:

1. Clustering is usually done by comparing each user's reports with other users. We aim to speed up this clustering by using latent variable matrix factorization and then clustering on much smaller latent variables.
2. Bootstrapping the clustering process by using non-task data to form an initial clustering before agent reporting begins.

2. Background and Related Work

In [?], Shnayder et al. prove that the Correlated Agreement multi-task peer prediction mechanism maximally incentivizes truthful reporting without a ground truth for homogeneous users. The CA mechanism can be generalized to non-homogeneous users but requires a prohibitive number of user reports; thus, [?] bounds expected regret by clustering users with similar signal-report distributions. This regret bound is proportional to the L1 difference between user-user ' Δ Matrices' and those of their respective clusters, where a Δ Matrix of users q, r is defined as:

$$D_{q,r}(i, j) = p(q \text{ reports } i, r \text{ reports } j) - p(q \text{ reports } i)p(r \text{ reports } j)$$

As far as we know, [?] is the only paper to specifically consider clustering methods for heterogeneous peer prediction. The paper differs from our work in that it assumes all of the relevant data exists and seeks to prove a theoretical upper bound on regret rather than to study empirical average regret in settings with sparse data. To demonstrate this, the authors restrict their dataset to a small subset of only those users for whom there is a full joint distribution between any two users, resulting in only 269 users.

Further, while the authors prove a theoretical bound in the case of empirically estimated ground truth, the tests they

run on real data rely on the existence of a known ground truth label. This allows them to calculate individual signal confusion matrices using methods developed in ? and leverage tensor decomposition techniques from ? and ? to compute average clusters directly. We seek to show that in the more general case of incomplete data with no ground truth labels, clustering instead on user-specific latent variables learned through matrix factorization can lead to similar computational savings as tensor decomposition while eliminating the need for labels. Further, we show that these computationally efficient clusterings are similarly effective with respect to incentivizing truthfulness as clustering on Δ matrices directly.

The use of the Correlated Agreement mechanism (?) requires that the Δ matrices are already known or can be empirically estimated from user reports. However, in practice, newly introduced users have contributed no reports and thereby cannot be clustered accurately. While nothing in the peer prediction literature addresses this cold start problem, it is well studied in the field of matrix factorization. The most similar method to ours, in that it incorporates non-task data, is (?). In this model, the author attempts to combine collaborative filtering and content-based recommendation systems by encoding item and user latent vectors as sums of vectors for the non-task features to which they correspond. Non-task features of an item or user in this model correspond to binary indicators of whether or not text labels describe the item or user. For example, in the setting of retail recommendations, features of a given item might be ‘pencil skirt’ and ‘blue’. Features of a user are text character descriptions drawn from user profiles. In our expansion on this model, we allow non-binary user features from the results of psychology and aptitude tests. We also allow for a portion of the user latent vectors to be unexplained by non-task user features and learned over time as tasks are completed.

3. Datasets

3.1. Task Data

The primary dataset we use is that of the *Good Judgment Global Forecasting Competition*, available at <https://dataverse.harvard.edu/dataverse.xhtml?alias=gjp>. There are many challenges inherent to this dataset. First, the idea of a task is not well-defined. Users enter predictions in continuous time, and the true probability of an event will change over time, such that each minute could be considered a different task (Note that if we expect users to only update their predictions in response to new information, these time-question tasks will be independent conditional on the signal). We choose to bucket these responses by week to mediate between the desire to allow for variation in true

probability over time while maintaining sufficient task response density. This results in 1,499,294 predictions by 2,052 users on 1,548 tasks. Secondly, the signal itself can be expected to have a high degree of noise, as the probability of any of these global events is not well-defined even in the absence of individual interpretation effects.

To study the effects of the clustering methodologies in a slightly less noisy realm, we additionally use a second dataset of user judgments on whether or not websites contained adult content. This dataset, termed *Adult*, was pulled from the Square Project at the University of Texas. The data is available at <https://github.com/ipeirotis/Get-Another-Label/tree/master/data>. We limit the dataset to users who gave at least one report of each signal, resulting in 65,144 reports by 389 users on 10,381 tasks.

Both datasets required significant preprocessing to account for sparsity of signal distributions and to compute the Δ matrices. Additionally, because Good Judgment reports were continuous rather than discrete, we bucketed the responses into five uniform bins, which we use to compute the matrices.

3.2. Non-Task Data

In addition to user predictions (task data), the *Good Judgment* dataset includes a variety of questions to assess the knowledge and ability of each user as well as their potential psychological biases. For instance, questions to assess knowledge might ask them about world events or to solve a math problem. Questions about character might ask them to state on a scale from 1-7 how just they believe the world is. We then use this non-task data to address the cold start problem discussed in section 2.

The psychological data contains both binary and ordinal features. A binary feature may be an answer to a ‘yes/no’ question or whether the user correctly answered a diagnostic math test question. An ordinal feature captures some continuous measure of that user’s character, *e.g.* user age or weekly workhours. To discretize ordinal features, we bucket each into a finite set of uniform bins similar to the Good Judgment report data. However, the number of bins is feature-dependent.

The dataset consists of 168 features for each of the 2,052 users. There are a significant number of missing answers in the dataset. For binary features, we made negative answers into -1 , missing answers into 0 , and positive answers into 1 . For ordinal features, we filled NaNs in with the average value of that feature. [Talk about data cleaning here?](#)

4. Model

Given a sparse matrix containing tasks, users, and users' reports on each task, the goal is to produce clusters which accurately describe users' reporting behavior. Because of the sparse matrix data, we focus on matrix factorization models to generate latent factor representations of users. That is, we assume that there is some latent variable representation of a user's forecasting behavior, u_i that interacts with related task-specific factors, t_j . These could be thought of, for example as the task's topic areas and the user's expertise or bias in those specific topic areas.

4.1. Good Judgment Matrix Factorization Model

After transforming the predictions to range from -1 to 1 instead of 0 to 1, we use the following model, shown in Figure ??:

$$p_{i,j} \sim \mathcal{N}(\tanh(\mu_{i_M} [u_i^T t_j + g + \mu_j] + \mu_{i_A}), \sigma^2)$$

Users tend to exhibit two forms of individualized biases in this space. The multiplicative bias, μ_{i_A} captures either risk aversion or overconfidence, such that they tend to predict values either closer to or farther from 0 (.5) than are warranted by their knowledge. The additive bias, μ_{i_M} captures both availability bias or normalcy bias which would result in user having a tendency to predict a specific level. The task bias μ_j captures the general prediction likelihood of the task absent user-specific factors, and the global bias g captures the overall mean of the user-task predictions. The tanh restricts the space of the predicted mean to the domain of the target value.

Then, we must infer $k + 2$ parameters for each user, $k + 1$ parameters for each task, and one global bias, where k is the length of the latent vector representation we choose for users and tasks. Only the predictions $p_{i,j}$ are known in this model.

For *Adult*, which does not have this probabilistic representation and has reports ranging from 1 to 4, we use only additive biases and remove the tanh.

4.2. Using Non-Task Data to Extend the Model for Cold Start Users

Clustering on the Δ matrix user representation requires a significant amount of task data before a user's appropriate cluster can be identified. Estimating *Good Judgment* matrix factorization model depends on an accurate estimation of various individual user parameters, requiring several predictions from a user before accurate predictions can be made. To mitigate these concerns, we use non-task psychological data to try to predict a users latent variables. This is easier than predicting a user's cluster as it leads to an intuitive distance metric for training, whereas trying to directly

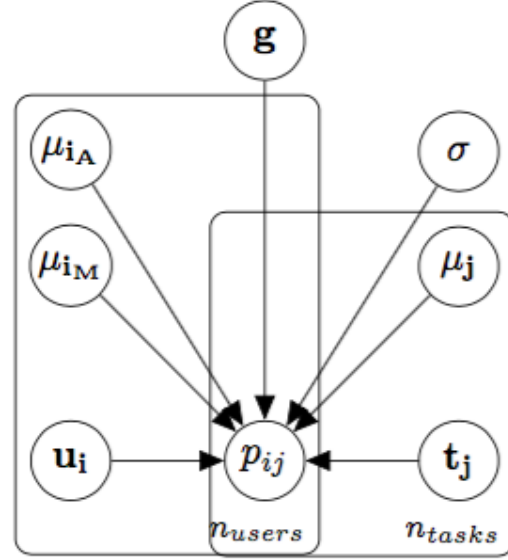


Figure 1. Model of prediction generation.

predict the user's cluster provides only a very hard binary classification task.

We focus only on the prediction of user latent vectors from non prediction data. Similar techniques could be used to solve the cold-start problem for tasks, however this would require data about tasks which come in before predictions. This is not part of the standard Good Judgement dataset and we did not explore feature engineering to create pre-prediction task data. An interesting direction for future work would be to consider natural language analysis techniques on text to create pre-prediction task data. For predicting the user latent variable models from psychological data, we try two separate techniques:

1. Predict latent variables from psychological data:

In our first model, we try predicting a user's latent variable representation directly from their psychological data. We use both a more interpretable simple linear regression model and a less interpretable fully connected neural network to perform the predictions. Then, it would be possible to make a prediction of a user's task report using a weighted combination of the predicted latent variables and the trained latent variables from matrix factorization.

2. Create user latent variable representation by summing feature latent variables

In this model, adapted from (?), a user's latent variable representation is created through a summation of various feature latent variables. If f is a set of features, then each feature has a latent vector \mathbf{L}_f . Then for user i with feature values u_{if} , the user's latent variable

representation is

$$u_i = \sum_{f=1}^{|f|} u_{if} \mathbf{L}_f + \mathbf{L}_{u_i}$$

where \mathbf{L}_{u_i} is a user-specific latent vector which starts at 0. Then, predictions for a user start by using latent vector $\sum_{f=1}^{|f|} u_{if} \mathbf{L}_f$ only as the user-specific latent variable starts at 0. As data on the user comes in, the user-specific latent vector \mathbf{L}_{u_i} moves away from 0. Predictions are then made as usual through

$$p_{i,j} = \tanh(\mu_{i_M} [u_i^T t_j + g + \mu_j] + \mu_{i_A})$$

5. Training

For learning the latent user parameters, we use mini-batch stochastic gradient descent to find the maximum likelihood estimate. We choose not to use a prior for regularization as we find that the training and validation performance do not generally diverge. This is likely because there is already enough noise in the training data for each parameter that it is difficult to overfit the model. We do not explicitly code the initialization of the model, so of course it is possible it finds a local rather than global minimum. However, the performance of the model is consistent across random starts.

Algorithm 1 Mini-Batch SGD Parameter Fitting

```

learning rate =  $\alpha$ 
for epoch in epochs do
  for batch in batches do
    Calculate the gradient of the MSE Loss for each parameter for all users and tasks in the batch
    Update the parameters:
    param  $\leftarrow$  param -  $\alpha * \frac{d}{dparam}(\text{MSE Loss})$ 
  end for
end for

```

In the above algorithm, parameters include user and task specific biases, the global bias, the latent vectors for tasks, and the latent vectors for specific users as well as user features. In (?), the user feature latent vectors were trained at the same time as user specific latent vectors. However, for reasons described below, we find that it makes more sense to train the latent variable model first only for user feature latent vectors. Then, after locking the user feature latent vectors and a corresponding set of task parameters: t , we then train a set of user specific latent vectors. Reports for new users are predicted as

$$p_{i,j} = \tanh(u_i^T t_j + g + \mu_j)$$

where u_i is formed as described in Section ?? . Reports for all other users are predicted as

$$p_{i,j} = \tanh(\mu_{i_M} [u_i^T t_j + g + \mu_j] + \mu_{i_A})$$

The global bias g and task-specific bias μ_j were allowed to vary throughout both training stages, although their values did not change much from stage 1 to 2.

There are two reasons that we train the user feature vectors before the individual user biases. First, while the model incorporating user feature biases is clearly more general than vanilla matrix factorization, it has more parameters and can be sensitive to the learning rate. In particular, we found that training all at once provided a worse RMSE on the predicted report values than vanilla matrix factorization, with the former having a RMSE of 0.33 and the latter having an RMSE of 0.28. This is most likely due to a too high learning rate, but the model already takes a significant time to converge at the high learning rate. If we train the user feature vectors first, and then the individual user latent vectors, the final RMSE was 0.24.

Additionally, the model above does not a unique global minimum. To see this, assume we are at some global minimum and there exists binary feature f . Then consider subtracting some k dimensional vector a_k from the feature latent vector \mathbf{L}_f and then adding the same vector to each user latent vector L_{u_i} with u_i having feature f . Then, $u_i = \sum_{f=1}^{|f|} u_{if} \mathbf{L}_f + \mathbf{L}_{u_i}$ is the same for each user and so we still have the same globally minimum error. Thus, we seek to train our data in such a way that the user feature latent vectors are informative as possible. One way to achieve this is to train only the user feature latent variables first and minimize loss, and then to train the individual user latent vectors after.

6. Methods

For the *Good Judgment* Dataset, we use the file `survey_fcsts.yr4.tab`, consisting of all user forecasts for the fourth year of the tournament. The file has 1,685,784 lines, each consisting of details about a specific user id and forecast. We restrict to only binary valued questions, for which a single probability of ‘yes’ is given as the answer (as opposed to, for example, “Who will be inaugurated as President of Russia in 2012?”, which had 3 possible answers). We restrict to only users who responded to at least 30 Individual Forecasting Problems (IFPs) over the course of the year. Given that users can (and ideally do) update their predictions for IFPs over time, any given time period could be thought of as a new IFP. Following the scoring procedure of the Good Judgment, we assume that if a user enters a prediction and then does not update it, this implies that they would like their existing prediction to be carried forward until they do update. We allow each week an IFP is active to be a different prediction task to balance temporal variability in task signals with our desire for higher task report density.

In computing the user latent vectors, we transform each prediction p from the range $[0, 1]$ to $[-1, 1]$ via the transformation

$$p = (p - 0.5) * 2$$

. This puts the predictions over the entire range of the tanh function and makes it so the the multiplicative user bias smoothly transforms data towards predictions of 1 or -1 (as opposed to outside the range). We train with stochastic gradient descent for 75 epochs with a learning rate of .2.

6.1. Clustering

As a metric of the clustering success, we use the average element-wise L1 distance between clustered Δ matrices and the original user-user Δ matrices for each pair of users. This requires binning the continuous predictions into discrete buckets. We choose to use 5 uniformly spaced bins from 0 to 1. We then compute the empirical joint distribution for any two users based on these discrete signals, as well as the empirical signal distributions of each individual. As mentioned in Section ??, the Δ matrices can be calculated from these as:

$$D_{q,r}(i, j) = p(q \text{ reports } i, r \text{ reports } j) - p(q \text{ reports } i)p(r \text{ reports } j)$$

To cluster, we use the scikit-learn implementation of K-means on all cores. The baseline we compare to is that of clustering on the Δ matrices directly, which is the method described in ? and should in general be better at guaranteeing similarity of Δ matrices than clustering on other factors. In regards to user coordinates for clustering, we append all of the Δ matrices between that user and all other users, leaving NaNs where not enough data exists to compute a distribution, and flatten this into a vector. When clustering, we iteratively impute coordinate means for NaNs such that they are not taken into consideration in the clustering.

The *Adult* Dataset has 92,720 lines, each of which has a report by a user id on a website. Ratings are discrete and ordinal, and we use them as given for both learning latent user variables and for computing Δ matrices.

6.2. Cold Start Problem

As described in Section ??, in the cold start problem we focus on predicting latent variables directly instead of on predicting clusters. In technique 1, we predict latent variables from labeled data coming from vanilla matrix factorization. So this takes as input the results of a complete matrix factorization, so that we have labeled data for each user of (user id, psychological data) \times (latent vector representation). The first technique we tried was to perform simple linear regression over the psychological features onto

Table 1. Distribution of L1 Distances Across Clusterings

	Δ (100)	K=2 (100)	Δ (500)	K=2 (500)
mean	.0168	.0215	.0145	.0192
std	.0243	.0290	.0218	.0258
min	0	0	0	0
25%	.0021	.0043	.0015	.0036
50%	.0084	.0108	.0071	.0099
75%	.0216	.0271	.0186	.0244
max	.7118	.6864	.6536	.6191

Table 2. Runtime (s) of KMeans algorithm on User Representation

Clusters	Δ	K=2	K=3	K=4
100	423	.57	.55	.65
150	625	.66	.66	.66
200	800	.66	.77	.87
250	959	.87	.88	.97
300	1187	.96	1.1	1.1
350	1198	1.1	1.2	1.2

the latent vector representations. In the second case, we trained a fully connected neural network. We tried 3 and 4 layers for the network, and varied the number of internal nodes between 20 and 80 for both layers. We found the best performance with a shallow network of just 3 layers and a middle layer of 20 nodes. Both techniques use out of the box PyTorch libraries.

For the second technique of using user feature dependent latent vectors, we implemented the model using PyTorch. The model was then trained using stochastic gradient descent, with a learning rate of 0.2. We took 75 iterations to train the matrix factorization unless otherwise noted.

7. Results

7.1. Sophie

While we achieved the lowest training and validation loss using 10 latent user factors (See Figure ??), the best L1 loss against the original Δ matrices was achieved by clustering on only 2 latent user factors (See Figure ??). Distribution results are given in Table ??.

However, clustering on the latent factors for all values of K was significantly faster than clustering directly on the Δ matrices, as the flattened concatenated matrices result in 51,300 features for each user. See Table ?? for details.

7.2. Cold Start Problem

Table ?? shows the performance using a Neural Network and simple linear regression in predicting the latent variables for new users from just psychological data. The

Table 3. Cold Start Predictions

Method	Mean Distance
Baseline (Mean Latent Vector)	0.67
Linear Regression	0.58
Neural Network	0.56

mean distance is calculated as the ℓ_2 distance from their final trained values after observing predictions and updating their latent vectors. Both are compared to just using a latent vector made up of coordinate-wise mean values for each new user. For each technique, there is some improvement over just using the mean latent vector, and we see some improvement over a random guess for new users. Ideally, this effect would be much larger; the extent to which this is possible using the psychological data for Good Judgment is unknown. Indeed, considering that there was no prior method of predicting new user's reports for Good Judgment, this provides some utility.

8. Discussion

The results show us that it is possible to use a significantly compressed representation of user heterogeneity and still arrive at a clustering that is still almost as successful at encouraging truthful reporting as the significantly more time- and memory-intensive task of clustering directly on the user-user signal distributions. Notably, the user representation we use does not explicitly encode the interactions between any two users' reporting behaviors yet performs nearly as well as the clustering with full knowledge of this data.

We note that the RMSE and L1 loss of the training model are lowest for larger numbers of latent factors, yet the clustering performs best when only two latent factors are used. One possible explanation is that this setup forces more of the user variability into the multiplicative and additive biases, which are more specifically relevant to the task of predicting a distribution of reports given signals. Supporting this interpretation, clustering the *Adult* dataset, which does not have the task structure to allow for modeling known psychological biases (specifically with respect to probability reporting) as its latent user biases, does not degrade as we increase the number of latent factors.

TODO INCLUDE CHART

9. Conclusion

- What happened?
- What next?

Quisque ullamcorper placerat ipsum. Cras nibh. Morbi vel

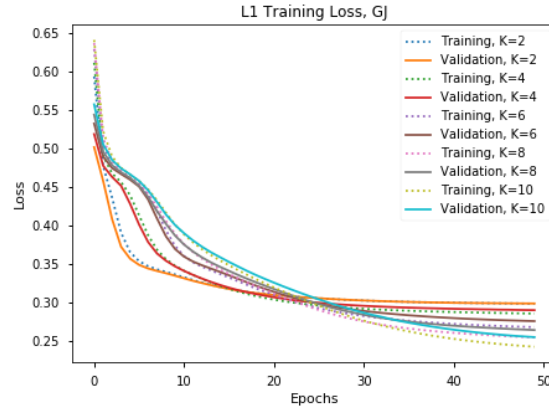


Figure 2. Training Results, L1 Loss.

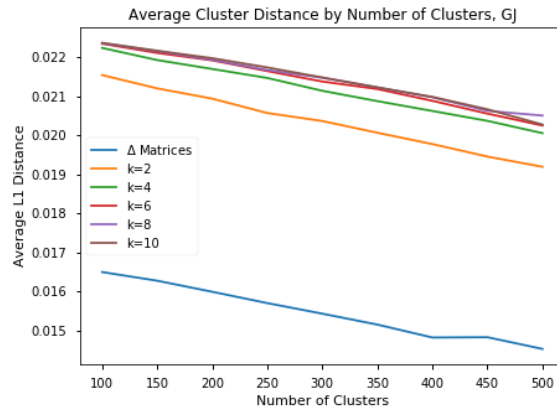


Figure 3. Average L1 Distance.

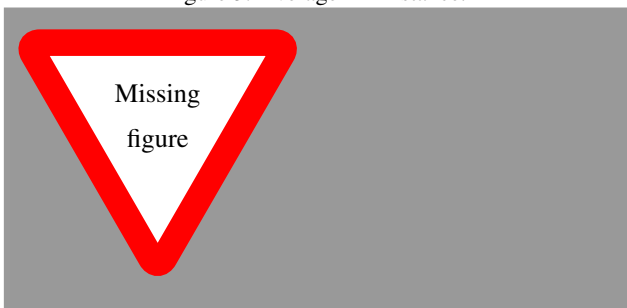


Figure 4. Visualizations of the internals of the system.

justo vitae lacus tincidunt ultrices. Lorem ipsum dolor sit amet, consectetur adipiscing elit. In hac habitasse platea dictumst. Integer tempus convallis augue. Etiam facilisis. Nunc elementum fermentum wisi. Aenean placerat. Ut imperdiet, enim sed gravida sollicitudin, felis odio placerat quam, ac pulvinar elit purus eget enim. Nunc vitae tortor. Proin tempus nibh sit amet nisl. Vivamus quis tortor vitae risus porta vehicula.

Fusce mauris. Vestibulum luctus nibh at lectus. Sed bibendum, nulla a faucibus semper, leo velit ultricies tellus, ac venenatis arcu wisi vel nisl. Vestibulum diam. Aliquam pellentesque, augue quis sagittis posuere, turpis lacus congue quam, in hendrerit risus eros eget felis. Maecenas eget erat in sapien mattis porttitor. Vestibulum porttitor. Nulla facilisi. Sed a turpis eu lacus commodo facilisis. Morbi fringilla, wisi in dignissim interdum, justo lectus sagittis dui, et vehicula libero dui cursus dui. Mauris tempor ligula sed lacus. Duis cursus enim ut augue. Cras ac magna. Cras nulla. Nulla egestas. Curabitur a leo. Quisque egestas wisi eget nunc. Nam feugiat lacus vel est. Curabitur consectetur.

Suspendisse vel felis. Ut lorem lorem, interdum eu, tincidunt sit amet, laoreet vitae, arcu. Aenean faucibus pede eu ante. Praesent enim elit, rutrum at, molestie non, nonummy vel, nisl. Ut lectus eros, malesuada sit amet, fermentum eu, sodales cursus, magna. Donec eu purus. Quisque vehicula, urna sed ultricies auctor, pede lorem egestas dui, et convallis elit erat sed nulla. Donec luctus. Curabitur et nunc. Aliquam dolor odio, commodo pretium, ultricies non, pharetra in, velit. Integer arcu est, nonummy in, fermentum faucibus, egestas vel, odio.

References

- Agarwal, Arpit, Mandal, Debmalya, Parkes, David C, and Shah, Nisarg. Peer prediction with heterogeneous users. 2017.
- Anandkumar, Animashree, Ge, Rong, Hsu, Daniel, Kakade, Sham M, and Telgarsky, Matus. Tensor decompositions for learning latent variable models. *Journal of Machine Learning Research*, 15:2773–2832, 2014.
- Dawid, Alexander Philip and Skene, Allan M. Maximum likelihood estimation of observer error-rates using the em algorithm. *Applied statistics*, pp. 20–28, 1979.
- Kula, Maciej. Metadata embeddings for user and item cold-start recommendations. *arXiv preprint arXiv:1507.08439*, 2015.
- Shnayder, Victor, Agarwal, Arpit, Frongillo, Rafael, and Parkes, David C. Informed truthfulness in multi-task peer prediction. In *Proceedings of the 2016 ACM Conference on Economics and Computation*, pp. 179–196. ACM, 2016.
- Zhang, Yuchen, Chen, Xi, Zhou, Dengyong, and Jordan, Michael I. Spectral methods meet em: A provably optimal algorithm for crowdsourcing. *The Journal of Machine Learning Research*, 17(1):3537–3580, 2016.