

Supplementary Material for LMLFM

Junjie Liang,^{1,2} Dongkuan Xu,² Yiwei Sun,^{1,3} Vasant Honavar^{1,2,3,4}

¹ Artificial Intelligence Research Laboratory, Pennsylvania State University

² College of Information Sciences and Technology, Pennsylvania State University

³ Department of Computer Science and Engineering, Pennsylvania State University

⁴ Institute of Computational and Data Sciences, Pennsylvania State University

{jul672,dux19,vhonavar}@ist.psu.edu, yus162@psu.edu

Derivation of LMLFM

In this section, we detail the process of solving LMLFM. Fig. outlines the hierarchical Bayesian model of LMLFM. The resulting generative model is as follows:

$$\begin{aligned} (y_{i0} | \mathbf{x}_{i0}, \boldsymbol{\theta}_i, \boldsymbol{\theta}_0, \alpha) &\sim N(y_{i0} | \hat{y}_{i0}, \alpha^{-1}) & (\alpha | \alpha_0, \beta_0) &\sim \text{Gamma}(\alpha | \alpha_0, \beta_0) \\ (\theta_{ik} | \mu_k^{\mathcal{I}}, b_k^{\mathcal{I}}) &\sim \text{Laplace}(\theta_{ik} | \mu_k^{\mathcal{I}}, b_k^{\mathcal{I}}) & (\theta_{ok} | \mu_k^{\mathcal{O}}, b_k^{\mathcal{O}}) &\sim \text{Laplace}(\theta_{ok} | \mu_k^{\mathcal{O}}, b_k^{\mathcal{O}}) \\ (\mu_k^{\mathcal{I}} | b_{\mu_0}) &\sim \text{Laplace}(\mu_k^{\mathcal{I}} | 0, b_{\mu_0}) & (\mu_k^{\mathcal{O}} | b_{\mu_0}) &\sim \text{Laplace}(\mu_k^{\mathcal{O}} | 0, b_{\mu_0}) \\ (b_k^{\mathcal{I}} | b_{b_0}) &\sim \text{Laplace}(b_k^{\mathcal{I}} | 0, b_{b_0}) & (b_k^{\mathcal{O}} | b_{b_0}) &\sim \text{Laplace}(b_k^{\mathcal{O}} | 0, b_{b_0}) \end{aligned}$$

we consider solving the Maximum A Posteriori (MAP) problem. Thus, the objective function is expressed as:

$$\Theta^* = \arg \max_{\Theta} \pi(\Theta | \mathbf{y}, X, \Theta_0) \quad (1)$$

where the model parameters and hyperpriors are represented by $\Theta = \{\alpha, \Theta, \mathbf{b}, \boldsymbol{\mu}\}$ and $\Theta_0 = \{\alpha_0, \beta_0, b_{b_0}, b_{\mu_0}\}$ respectively. In the following, we only provide the derivation w.r.t. the individual parameters, i.e., $\Theta^{\mathcal{I}} = \{\alpha, \Theta^{\mathcal{I}}, \boldsymbol{\mu}^{\mathcal{I}}, \mathbf{b}^{\mathcal{I}}\}$. To keep the notation light, we omit the superscript \mathcal{I} till the end of this section.

Update of $\Theta^{\mathcal{I}}$: For each model parameter $\theta_i \in \Theta$, the prediction is a linear combination of two functions $g(i)$ and $h(i)$ that are independent of the value of θ_i :

$$\hat{\mathbf{y}}_i = g(i) + h(i)\theta_i \quad (2)$$

with

$$g(i) = \text{diag}(X_i \cdot \Theta_i^{\mathcal{O}\top}) \quad h(i) = X_i + \Theta_i^{\mathcal{O}}$$

where $\Theta_i^{\mathcal{O}}$ is the matrix of latent factors constructed by the observations associated to i . Thus we have $\pi(\theta_i) = \pi(\mathbf{y}_i | X_i, \theta_i, \alpha) \cdot \pi(\theta_i | \mu_k, b_k)$, which is given by:

$$\pi(\theta_i) \propto \exp\left\{-\frac{\alpha}{2} \|\mathbf{y}_i - \hat{\mathbf{y}}_i\|_2^2\right\} \cdot \prod_{k=1}^p \exp\left\{-\frac{|\theta_{ik} - \mu_k|}{b_k}\right\}$$

The log of which is:

$$\ell(\theta_i) = -\frac{\alpha}{2} \|\mathbf{y}_i - \hat{\mathbf{y}}_i\|_2^2 - \sum_{k=1}^p \frac{|\theta_{ik} - \mu_k|}{b_k} + \text{const.} \quad (3)$$

Copyright © 2020, Association for the Advancement of Artificial Intelligence (www.aaai.org). All rights reserved.

Let's define the gradient of $|\theta_{ik} - \mu_k|$ by e . Then following the sub-gradient equations (see e.g., (Bertsekas 1999)), we have $e = \text{sgn}(\theta_{ik} - \mu_k)$ if $\theta_{ik} \neq \mu_k$ and $e \in [-1, 1]$ otherwise. For cases where $|\theta_{ik} - \mu_k|$ is differentiable, the optimal θ_{ik}^* can be derived by simply setting the gradient of (3) to 0. Otherwise, $\theta_{ik} = \mu_k$. Consequently, we have the following results:

$$\theta_{ik}^* = (\mathbf{h}_{ik}^{\top} \mathbf{h}_{ik})^{-1} \left(\mathbf{h}_{ik}^{\top} \mathbf{h}_{ik} \mu_k + \text{sgn}(\mathbf{r}_{ik}) (|\mathbf{r}_{ik}| - 1/\alpha b_k^{\mathcal{I}})_+ \right) \quad (4)$$

where $(\cdot)_+$ is the ReLU function; $\mathbf{r}_{ik} = \mathbf{h}_{ik}^{\top}(\mathbf{y}_i - g(i) - \sum_{q \in 1:p \setminus k} \theta_{iq} \mathbf{h}_{iq} - \mathbf{h}_{ik} \mu_k)$, with $1 : p \setminus k$ denoting the set of integers ranging from 1 to p excluding k and \mathbf{h}_{ik} is the k -th column of $h(i)$.

Update of α : $\pi(\alpha) = \pi(\mathbf{y} | X, \Theta, \alpha) \cdot \pi(\alpha | \alpha_0, \beta_0)$. The full conditional posterior of α can be obtained by:

$$\begin{aligned} \pi(\alpha) &\propto \prod_{j=1}^{|\mathbf{y}|} \alpha^{1/2} \exp\left\{-\frac{\alpha}{2} (y_j - \hat{y}_j)^2\right\} \cdot \alpha^{\alpha_0-1} \exp\{-\alpha \beta_0\} \\ &\propto \alpha^{\alpha_0 + |\mathbf{y}|/2 - 1} \exp\left\{-\alpha \left(\beta_0 + \|\mathbf{y} - \hat{\mathbf{y}}\|_2^2 / 2\right)\right\} \\ &\propto \text{Gamma}\left(\alpha_0 + |\mathbf{y}|/2, \beta_0 + \|\mathbf{y} - \hat{\mathbf{y}}\|_2^2 / 2\right) \end{aligned}$$

the mode of $\pi(\alpha)$ is then accessed by $(\beta_0 + \|\mathbf{y} - \hat{\mathbf{y}}\|_2^2 / 2)^{-1} (\alpha_0 + |\mathbf{y}|/2 - 1)$.

Update of $\boldsymbol{\mu}$: Since the components in $\boldsymbol{\mu}$ are i.i.d., we only consider computing the full conditional posterior of μ_k ($k = 1, \dots, p$). $\pi(\mu_k) = \pi(\boldsymbol{\theta}_{\cdot k} | \mu_k, b_k) \cdot \pi(\mu_k | 0, b_{\mu_0})$. The full conditional posterior of μ_k can be obtained by:

$$\pi(\mu_k) \propto \prod_{i=1}^n \exp\left\{-\frac{|\theta_{ik} - \mu_k|}{b_k}\right\} \cdot \exp\left\{-\frac{|\mu_k|}{b_{\mu_0}}\right\}$$

The problem of finding the optimal μ_k can be converted to finding the weighted median of the vector $\boldsymbol{\theta}_{\cdot k} \cup \{0\}$ with the weights $\{b_k\}_{i=1}^p \cup \{b_{\mu_0}\}$, which, as discussed in (Gurwitz 1990), can be solved in linear time.

Update of \mathbf{b} : Since the components in \mathbf{b} are i.i.d., we only consider computing the full conditional posterior of b_k ($k = 1, \dots, p$). $\pi(b_k) = \pi(\boldsymbol{\theta}_{\cdot k} | \mu_k, b_k) \cdot \pi(b_k | 0, b_{b_0})$. The

Algorithm 1 LMLFM

- 1: **Input:** Training set $S = \{X, \mathbf{y}\}$, hyperpriors Θ_0 , convergence criteria T (See Sec.).
 - 2: **Output:** Θ .
 - 3: Initialize the parameters $\{b, \alpha, \mu, \Theta\}$
 - 4: **repeat**
 - 5: Update the parameters in the following order:
 $\{\Theta, \alpha, \mu^{\mathcal{I}}, \mu^{\mathcal{O}}, b^{\mathcal{I}}, b^{\mathcal{O}}\}$
 - 6: **until** convergence
-

full conditional posterior of b_k can be obtained by:

$$\pi(b_k |) \propto \prod_{i=1}^n b_k^{-1} \exp\left\{-\frac{|\theta_{ik} - \mu_k|}{b_k}\right\} \cdot \exp\left\{-\frac{|b_k|}{b_{b_0}}\right\}$$

Assume $b_k > 0$, by setting $d \log \pi(b_k |) / db_k = 0$, we arrive at:

$$\frac{b_k^2}{b_{b_0}} + nb_k - \|\boldsymbol{\theta}_{\cdot k} - \mu_k\|_1 = 0$$

This is a simple quadratic function. Combining the fact that $b_k > 0$, the root of is found by:

$$b_k^* = 2b_{b_0} \left(\sqrt{n^2 + \frac{4}{b_{b_0}} \|\boldsymbol{\theta}_{\cdot k} - \mu_k\|_1} - n \right) \quad (5)$$

We can easily verify that $b_k > 0$ only when $\|\boldsymbol{\theta}_{\cdot k} - \mu_k\|_1 > 0$. For the case where $\|\boldsymbol{\theta}_{\cdot k} - \mu_k\|_1 = 0$, we have $\theta_{ik} = \mu_k$ for all $i = 1, \dots, n$. This is equivalent to rejecting the randomness of $\boldsymbol{\theta}_{\cdot k}$; thus, if $\|\boldsymbol{\theta}_{\cdot k} - \mu_k\|_1 = 0$, we can draw two conclusions: i) variable k is subject to fixed effect instead of random effects; ii) $\boldsymbol{\theta}_{\cdot k} \sim \text{Laplace}(\mu_k, 0)$, thus we get $b_k = 0$.

The pseudo-code of LMLFM is given in Algorithm 1.

Complexity Analysis.

The computation of Θ can be accelerated by pre-computing and caching $\hat{\mathbf{y}}$. Then, $\hat{\mathbf{y}}$ can be maintained up-to-date using:

$$g(i) + \sum_{q \in 1:p \setminus k} \theta_{iq} \mathbf{h}_{iq} = \hat{\mathbf{y}}_i - \theta_{ik}^{(t)} \mathbf{h}_{ik}; \quad \hat{\mathbf{y}}_i = \hat{\mathbf{y}}_i + \theta_{ik}^{(t+1)} \mathbf{h}_{ik}$$

where $\theta_{ik}^{(t)}$ is the value of θ_{ik} for the t -th iteration. Thus, updating Θ takes $O(|S|)$ time. Next, computing α requires $\hat{\mathbf{y}}$, which is already cached; hence, complexity for updating α is $O(|\mathbf{y}|)$. The time complexity of updating μ, b is $O(p(n+m))$. Therefore, the overall computational complexity for one complete iteration is $O(|S|)$, which is strictly linear in the size of the training data. It is worth noting that the time complexity of FM is $O(k|S|)$ (Rendle 2012), where k denotes the number of latent factors. The space complexity of LMLFM is $O(|\mathbf{y}|)$, identical to that of FM (Rendle 2012). We conclude that our algorithm is more efficient than FM.

Theoretical Analysis

In this section, we will prove two important properties in LMLFM: *Ascent property* and *Convergence*. Let \mathbb{Z}^+ denotes the set of all positive integers. To keep the notation light,

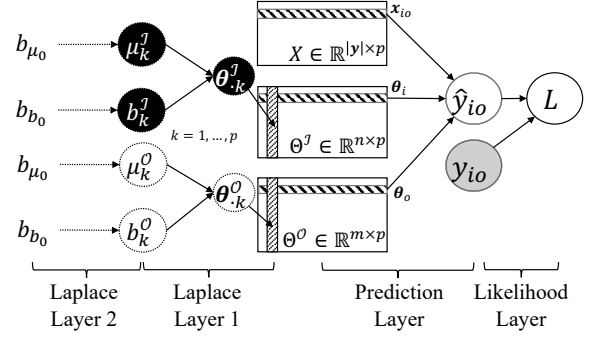


Figure 1: The hierarchical Bayesian structure of LMLFM. Laplace layer 1 is designed to select predictive latent factors; Laplace layer 2 is designed to identify random effects.

we overload the symbol θ_j to refer to the j -th component of the vectorized full parameter set $\Theta = \{\theta_j\}_{j=1}^d$. The remaining set of parameters excluded θ_j is denoted by Θ_{-j} . We denote by $\pi(\theta_j |)$ the full conditional posterior of θ_j , i.e., $\pi(\theta_j | \mathbf{y}, X, \Theta_0, \Theta_{-j})$. Putting these together, the joint posterior density at the t -th iteration in (1) can be expanded as:

$$\pi(\Theta^{(t)} | \mathbf{y}, X, \Theta_0) = \pi(\theta_j^{(t)} |) \cdot \pi(\Theta_{-j}^{(t)} | \mathbf{y}, X, \Theta_0) \quad (6)$$

Proposition 1. Ascent property. $\pi(\Theta^{(t+1)}) \geq \pi(\Theta^{(t)})$ holds for all iteration $t \in \mathbb{Z}^+$.

Proof. Without loss of generality, we assume that Θ is updated with the following order $\theta_1, \dots, \theta_d$, where θ_j is updated by the mode of $\pi(\theta_j |)$, $j = 1, \dots, d$. Therefore, when $j = 1$, we must have $\pi(\theta_1^{(t+1)} |) \geq \pi(\theta_1^{(t)} |)$. Thus, followed by (6), we have $\pi(\theta_1^{(t+1)}, \theta_2^{(t)}, \dots, \theta_d^{(t)} |) \geq \pi(\theta_1^{(t)}, \theta_2^{(t)}, \dots, \theta_d^{(t)} |)$. Similarly, for $j = 2, \dots, d$, $\theta_j^{(t+1)} = \arg \max_{\theta_j} \pi(\theta_j |)$ with other parameters fixed. Consequently, we have $\pi(\theta_1^{(t+1)}, \dots, \theta_j^{(t+1)}, \theta_{j+1}^{(t)}, \dots, \theta_d^{(t)} |) \geq \pi(\theta_1^{(t+1)}, \dots, \theta_{j-1}^{(t+1)}, \theta_j^{(t)}, \theta_{j+1}^{(t)}, \dots, \theta_d^{(t)} |)$. At $j = d$, we achieve $\pi(\Theta^{(t+1)}) \geq \pi(\Theta^{(t)})$. \square

Proposition 2. Convergence. If $\pi(\Theta^{(t)})$ is bounded above, there exists an iteration $t \in \mathbb{Z}^+$, such that $\forall i \in \mathbb{Z}^+, |\pi(\Theta^{(t+i)}) - \pi(\Theta^{(t)})| < \epsilon$ holds for $\epsilon > 0$.

Proof. According to Proposition 1, we only need to prove $\pi(\Theta^{(t+i)}) - \pi(\Theta^{(t)}) < \epsilon$. Let's assume that $\forall t \in \mathbb{Z}^+$, such that $\exists i \in \mathbb{Z}^+, \pi(\Theta^{(t+i)}) - \pi(\Theta^{(t)}) \geq \epsilon$. Then we can always find a monotonically increasing sequence $\{a_k\}_{k=0}^\infty$ with $a_0 = 0$, such that $\pi(\Theta^{(t+a_k)}) > \pi(\Theta^{(t+a_{k-1})}) > \dots > \pi(\Theta^{(t+a_0)})$. Let's further define ϵ_k as the lower bound of $\pi(\Theta^{(t+a_k)}) - \pi(\Theta^{(t+a_{k-1})})$. Thus, $\sum_{j=1}^k \pi(\Theta^{(t+a_j)}) -$

$\pi(\Theta^{(t+a_{j-1})}) = \pi(\Theta^{(t+a_k)}) - \pi(\Theta^{(t)}) \geq \sum_{j=1}^k \epsilon_j$. Then we have $\pi(\Theta^{(t+a_k)}) \geq \pi(\Theta^{(t)}) + \sum_{j=1}^k \epsilon_j$. Notice that $\forall k, \epsilon_k > 0$, we then have $\lim_{k \rightarrow \infty} \pi(\Theta^{(t+a_k)}) \geq \pi(\Theta^{(t)}) + \sum_{j=1}^{\infty} \epsilon_j \geq \sum_{j=1}^{\infty} \epsilon_j \rightarrow \infty$. However, the density function $\pi(\Theta)$ is assumed bounded, thus leads to a contradiction. \square

Proposition 1 establishes the convergence of the joint posterior density $\pi(\Theta)$. It also allows us to set the stopping criteria of our model: checking either the number of iterations exceeds a predefined threshold or the improvement of joint posterior density is below a preset threshold. However, because computing $\pi(\Theta)$ is intractable, we instead monitor the full joint density $\pi(\mathbf{y}, X, \Theta, \Theta_0)$.

Correlation Estimation

Longitudinal Correlation (LC) Let's denote o, j as two different observations. The LC is computed as follows:

$$\begin{aligned} \text{cov}(y_{io}, y_{ij}) &= \mathbb{E}[(\mathbf{x}_{io}^\top(\boldsymbol{\theta}_i + \boldsymbol{\theta}_o) + \boldsymbol{\theta}_i^\top \boldsymbol{\theta}_o)(\mathbf{x}_{ij}^\top(\boldsymbol{\theta}_i + \boldsymbol{\theta}_j) + \boldsymbol{\theta}_i^\top \boldsymbol{\theta}_j)^\top] \\ &\quad - \mathbb{E}[(\mathbf{x}_{io}^\top(\boldsymbol{\theta}_i + \boldsymbol{\theta}_o) + \boldsymbol{\theta}_i^\top \boldsymbol{\theta}_o)] \cdot \mathbb{E}[(\mathbf{x}_{ij}^\top(\boldsymbol{\theta}_i + \boldsymbol{\theta}_j) + \boldsymbol{\theta}_i^\top \boldsymbol{\theta}_j)] \\ &= 2(\mathbf{x}_{io} + \boldsymbol{\mu}^\mathcal{O})^\top \llbracket \mathbf{b}^\mathcal{I} \rrbracket^2 (\mathbf{x}_{ij} + \boldsymbol{\mu}^\mathcal{O}) \end{aligned}$$

Cluster Correlation (CC) Let's denote i, v as two different individuals. The CC is computed as follows:

$$\begin{aligned} \text{cov}(y_{io}, y_{vo}) &= \mathbb{E}[(\mathbf{x}_{io}^\top(\boldsymbol{\theta}_i + \boldsymbol{\theta}_o) + \boldsymbol{\theta}_i^\top \boldsymbol{\theta}_o)(\mathbf{x}_{vo}^\top(\boldsymbol{\theta}_v + \boldsymbol{\theta}_o) + \boldsymbol{\theta}_v^\top \boldsymbol{\theta}_o)^\top] \\ &\quad - \mathbb{E}[(\mathbf{x}_{io}^\top(\boldsymbol{\theta}_i + \boldsymbol{\theta}_o) + \boldsymbol{\theta}_i^\top \boldsymbol{\theta}_o)] \cdot \mathbb{E}[(\mathbf{x}_{vo}^\top(\boldsymbol{\theta}_v + \boldsymbol{\theta}_o) + \boldsymbol{\theta}_v^\top \boldsymbol{\theta}_o)] \\ &= 2(\mathbf{x}_{io} + \boldsymbol{\mu}^\mathcal{I})^\top \llbracket \mathbf{b}^\mathcal{O} \rrbracket^2 (\mathbf{x}_{vo} + \boldsymbol{\mu}^\mathcal{I}) \end{aligned}$$

Multi-level Correlation Finally, the multi-level correlation is computed as follows:

$$\begin{aligned} \text{cov}(y_{io}, y_{io}) &= \mathbb{E}[(\mathbf{x}_{io}^\top(\boldsymbol{\theta}_i + \boldsymbol{\theta}_o) + \boldsymbol{\theta}_i^\top \boldsymbol{\theta}_o)(\mathbf{x}_{io}^\top(\boldsymbol{\theta}_i + \boldsymbol{\theta}_o) + \boldsymbol{\theta}_i^\top \boldsymbol{\theta}_o)^\top] \\ &\quad - \mathbb{E}[(\mathbf{x}_{io}^\top(\boldsymbol{\theta}_i + \boldsymbol{\theta}_o) + \boldsymbol{\theta}_i^\top \boldsymbol{\theta}_o)] \cdot \mathbb{E}[(\mathbf{x}_{io}^\top(\boldsymbol{\theta}_i + \boldsymbol{\theta}_o) + \boldsymbol{\theta}_i^\top \boldsymbol{\theta}_o)] \\ &= 2\mathbf{x}_{io}^\top (\llbracket \mathbf{b}^\mathcal{I} \rrbracket^2 + \llbracket \mathbf{b}^\mathcal{O} \rrbracket^2) \mathbf{x}_{io} + 4 \cdot \text{tr}(\llbracket \mathbf{b}^\mathcal{I} \rrbracket^2 \cdot \llbracket \mathbf{b}^\mathcal{O} \rrbracket^2) \end{aligned}$$

where $\text{tr}(A)$ is the trace of matrix A .

Parametrization

Initialization. Previous study has pointed out that ICM-based techniques are particularly sensitive to the choice of initializations (Szeliski et al. 2008). We find in our experiments that LMLFM is rather robust to the choice of parameter initialization. The set of parameters that need initialization are $\{\mathbf{b}, \alpha, \boldsymbol{\mu}, \Theta\}$. We simply let $\mathbf{b} = \mathbf{1}$, $\boldsymbol{\mu} = \mathbf{0}$ and $\Theta = \mathbf{0}$. The precision α tends to have stronger impact on the model performance. Our experiments show that setting α with the form $\tau [\text{var}(\mathbf{y})]^{-1}$ (where τ is a hyper-parameter controlling the value of precision) provides satisfactory results.¹

¹ τ is the only tuning hyper-parameter in LMLFM.

Hyperprior settings. Due to the high number of explanatory variables, the impact of α_0, β_0 are negligible; therefore, we chose trivial values $\alpha_0 = \beta_0$. We also set $b_{\mu_0} = b_{b_0} = 1$. We claim convergence if the change of normalized log posterior density value of two consecutive iterations is less than 10^{-4} or the number of iterations exceed 10.

Experimental Protocol

The observations in the generated data are assigned to disjoint training and test sets while simultaneously ensuring that for any given individual, no observation that is included in the training data has a time stamp that is later than any observation that is included in the test data. This ensures that while making predictions no information from the future is used. Specifically, we use the following procedure: i) Randomly split the complete data set into two subsets, i.e., $S^{(a)}$ (70%) and $S^{(b)}$ (30%); ii) For each individual $i \in S^{(a)}$, we find out the observation with the latest time stamp and denote it as t_i ; iii) For each individual $i \in S^{(b)}$, we split the observations associated to i into two half. Specifically, let $S_i^{(b_1)}$ denotes the subset of observations with time stamp less than t_i and $S_i^{(b_2)}$ denoting the rest; iv) The training and test data is given by $S^{(a)} \cup S^{(b_1)}$ and $S^{(b_2)}$ respectively.

Most of the implementations of our baselines are publicly available. we use the LASSO code from (Pedregosa et al. 2011). For M-LMM, LMMLASSO, GLMMLASSO and rPQL, we utilize the `lmer4`, `lmmlasso`, `glmmlasso` and `rpql` package, respectively from CRAN.² We implemented MLLASSO and LMLFM in Python. We use the implementation of Random forest in (Pedregosa et al. 2011), FM in (Mikhail Trofimov 2016) and the Penalized GEE (PGEE) in the `PGEE` package in CRAN. All experiments are conducted on a desktop computer with i7-7700K CPU, 16G RAM and GeForce GTX 1060 graphics card. The hyperparameters of all methods are tuned to optimize their performance using 5-fold cross validation. We report performance statistics obtained from 100 independent runs. Evaluation scores are computed using only the test data.

Simulated data. The outcomes \mathbf{y} are drawn independently from $N(y_{io} | \mathbf{x}_{io}^\top(\boldsymbol{\beta} + \boldsymbol{\theta}_i + \boldsymbol{\theta}_o) + \boldsymbol{\theta}_i^\top \boldsymbol{\theta}_o, 1)$, with $\boldsymbol{\beta} = \{1, 2, 3, -1, -2, -3, 7, 10, 0, \dots, 0\}$. Each component $\theta_{zk} \in \Theta$ follows $Laplace(\theta_{zk} | 0, b_k)$, where $b_k^\mathcal{I}$ for $k = 1, \dots, 10$ and $b_k^\mathcal{O}$ for $k = 5, \dots, 15$ are drawn from $U(0, 1)$. For other k , we set $b_k^\mathcal{I} = b_k^\mathcal{O} = 0$. Thus there are 15 and $p - 15$ relevant features with random effects and fixed effects respectively. Absence of LC or CC is simulated by manually setting $\{\Theta^\mathcal{O}, \mathbf{b}^\mathcal{O}\}$ or $\{\Theta^\mathcal{I}, \mathbf{b}^\mathcal{I}\}$ to zero, respectively.

References

- Bertsekas, D. P. 1999. *Nonlinear programming*. Athena scientific Belmont.
- Gurwitz, C. 1990. Weighted median algorithms for l1 approximation. *BIT* 30(2):301–310.

²<https://cran.r-project.org/>

Mikhail Trofimov, A. N. 2016. tffm: Tensorflow implementation of an arbitrary order factorization machine. <https://github.com/geffy/tffm>.

Pedregosa, F.; Varoquaux, G.; Gramfort, A.; Michel, V.; Thirion, B.; Grisel, O.; Blondel, M.; Prettenhofer, P.; Weiss, R.; Dubourg, V.; et al. 2011. Scikit-learn: Machine learning in python. *Journal of machine learning research* 12(Oct):2825–2830.

Rendle, S. 2012. Factorization machines with libfm. *ACM Transactions on Intelligent Systems and Technology (TIST)* 3(3):57.

Szeliski, R.; Zabih, R.; Scharstein, D.; Veksler, O.; Kolmogorov, V.; Agarwala, A.; Tappen, M.; and Rother, C. 2008. A comparative study of energy minimization methods for markov random fields with smoothness-based priors. *IEEE transactions on pattern analysis and machine intelligence* 30(6):1068–1080.

Depression Analysis on SWAN Data		LASSO	FM	RF	MLLASSO
Rank	LMLEM				
1	NervousPast2Wks.a.lot (4.91)	HospitalStayPastYear.yes (7.33)	OtherEventUpset.not.upset (-0.05)	NervousPast2Wks.a.lot (0.13)	Race.white (10.27)
2	HospitalStayPastYear.yes (2.89)	HospitalStayPastYear.no (7.01)	STATUS.hysterectomy/both ovaries removed (-0.05)	MoodChangePast2Wks.never (0.08)	Race.black.american (9.27)
3	HospitalStayPastYear.no (2.82)	Stroke.no (6.9)	NervousPast2Wks.a.lot (0.05)	Stroke.no (0.04)	HadD&C.yes (7.04)
4	MoodChangePast2Wks.never (-2.36)	Stroke.yes (5.66)	Race.hispanic (0.05)	Height (0.04)	ViolentEventUpset.no (6.71)
5	Income.very.low (2.27)	NervousPast2Wks.a.lot (3.88)	Income.very.low (0.05)	DHAS (0.03)	QuitJobUpset.no (6.7)
6	Race.hispanic (2.26)	MoodChangePast2Wks.a.lot (3.2)	STATUS.unknown.due.to.hysterectomy (-0.05)	Testosterone (0.03)	QuitJobUpset.not.upset (6.68)
7	Stroke.no (1.77)	NervousPast2Wks.never (-2.72)	MoodChangePast2Wks.a.lot (0.04)	SHBG (0.03)	PartnerUnempUpset.no (6.63)
8	TalkLastYr.a.lot (1.69)	Height (-2.29)	STATUS.pregnant/breastfeeding (0.04)	FSH (0.03)	HipSize (6.51)
9	TalkLastYr.few (1.32)	Income.very.low (2.26)	HotFlashesPast2Wks.a.lot (-0.04)	Estradiol (0.03)	EndedRelationUpset.no (6.46)
10	WorsenRelationUpset.yes (1.25)	Migraines.yes (2.16)	STATUS.post-menopausal (-0.04)	HipBoneDensity (0.03)	HadD&C.no (6.36)
General Happiness Analysis on GSS Data		LASSO	FM	RF	MLLASSO
Rank	LMLEM				
1	MarriageHappiness.no (-0.53)	Financial.satisfied (0.21)	FirstLanguage.Tagalog (-1.53)	WomensHealthEndangered.yes (0.08)	Interviews.Language.English (2.47)
2	LifeExciting.dull (-0.45)	LifeExciting.dull (-0.2)	LanguageAtHome.Vietnamese (-0.99)	HaveLifeAfterDeath.yes (0.07)	MemberIn.NationalityGp.no (2.38)
3	JobOrHousework.very.dissatisfied (-0.21)	Financial.more.or.less.satisfied (0.19)	FirstLanguage.KOREAN (0.97)	CourtsDealWithCriminals.not.harse.enough (0.06)	MemberIn.SchoolFraternity.no (2.37)
4	Health.poor (-0.15)	Oppo.RaceInNeighborhood.no (0.17)	2ndCountrySpouseOrigin.WestIndies (-0.85)	LifeExciting.dull (0.03)	MemberIn.PoliticalClub.no (2.37)
5	MarriageHappiness.very.happy (0.11)	LifeExciting.exciting (0.17)	FirstLanguage.Hindu (-0.84)	Financial.not.at.all.satisfied (0.03)	MemberIn.FarmOrg.no (2.37)
6	Unemployed (-0.1)	MarriageHappiness.very.happy (0.17)	VoteIn1972Election.dont.known (-0.83)	EverUnempInLast10Yrs.no (0.03)	MemberIn.ArGp.no (2.35)
7	MarriageHappiness.pretty.happy (0.09)	Oppo.RaceInNeighborhood.yes (0.16)	3rdCountrySpouseOrigin.Philippines (0.83)	ScenX-RatedMovieLastYr.no (0.02)	MemberIn.ServiceGp.no (2.34)
8	JobOrHousework.dissatisfied (-0.08)	FinancialChange.better (0.16)	FirstLanguage.Norwegian (-0.79)	MarriageHappiness.very.happy (0.02)	MemberIn.YouthGp.no (2.33)
9	Financial.not.satisfied (-0.08)	Financial.stayed.same (0.13)	3rdCountryYouOrigin.WestIndies (0.78)	Oppo.RaceInNeighborhood.yes (0.01)	MemberIn.FraternalGp.no (2.33)
10	SocialClass.lower.class (-0.08)	PovertyStatus.unknown (-0.11)	2ndCountrySpouseOrigin.arabic (0.72)	PovertyStatus.unknown (0.01)	MemberIn.VeteranGp.no (2.33)

* D&C: Dilation and curettage; DHAS: Dehydroepiandrosterone sulfate; FSH: Follicle-stimulating hormone; SHBG: Sex hormone-binding globulin.

Table 1: Top 10 most relevant variable selected by different models on the real-life data sets. The feature importance (or PAE) is specified in the parenthesis.