

Federated Binary Matrix Factorization using Proximal Optimization

Sebastian Dalleiger¹ Jilles Vreeken² Michael Kamp³

¹KTH Royal Institute of Technology

²CISPA Helmholtz Center for Information Security

³Institute for AI in Medicine, UK Essen and Ruhr University Bochum, Monash University
sdall@kth.se jv@cispa.de michael.kamp@uk-essen.de

Abstract

Identifying informative components in binary data is an essential task in many application areas, including life sciences, social sciences, and recommendation systems. Boolean matrix factorization (BMF) is a family of methods that performs this task by factorizing the data into dense factor matrices. In real-world settings, the data is often distributed across stakeholders and required to stay private, prohibiting the straightforward application of BMF. To adapt BMF to this context, we approach the problem from a federated-learning perspective, building on a state-of-the-art continuous binary matrix factorization relaxation to BMF that enables efficient gradient-based optimization. Our approach only needs to share the relaxed component matrices, which are aggregated centrally using a proximal operator that regularizes for binary outcomes. We show the convergence of our federated proximal gradient descent algorithm and provide differential privacy guarantees. Our extensive empirical evaluation shows that our algorithm outperforms, in quality and efficacy, federation schemes of state-of-the-art BMF methods on a diverse set of real-world and synthetic data.

1 Introduction

Discovering patterns and dependencies in distributed binary data sources is a common problem in many applications, such as cancer genomics (Liang, Zhu, and Lu 2020), recommender systems (Ignatov et al. 2014), and neuroscience (Haddad et al. 2018). Data is often distributed horizontally (i.e., the rows of the data matrix are split across hosts) and may not be pooled. For example, biopsies are performed in different hospitals, with each location measuring the expression of a common set of genes. Although there exists an explicit interest in analyzing this data jointly, privacy regulations mandate that these measurements may not be shared, thereby limiting the applicability of traditional centralized methods.

Federated learning (McMahan et al. 2017) enables learning from distributed datasets *without* disclosing sensitive data. Existing methods for federated non-negative matrix factorization (Li et al. 2021) are specific to real-valued data, and similar to non-federated Non-negative Matrix Factorization (NMF) (Paatero and Tapper 1994; Lee and Seung 1999, 2000), singular value decomposition (Golub and Loan 1996),

and principal component analysis (Golub and Loan 1996), do not achieve interpretable results for binary data (Miettinen et al. 2008; Dalleiger and Vreeken 2022).

Boolean Matrix Factorization (BMF) alleviates this problem by approximating a *centralized* Boolean target matrix $A \in \{0, 1\}^{n \times m}$ by the Boolean product

$$A \approx [U \circ V]_{ij} = \bigvee_{l \in [k]} U_{il} \wedge V_{lj}$$

of two low-rank Boolean factor matrices (Miettinen et al. 2008), $U \in \{0, 1\}^{n \times k}$ (*feature matrix*) and $V \in \{0, 1\}^{k \times m}$ (*coefficient matrix*).

Although there are myriad heuristics to approximate this NP-hard problem, doing so for *distributed data* without sharing private information remains an open problem. Even though we could approach distributed binary data with standard federated learning techniques, e.g. aggregating locally-obtained BMF results into a shared matrix, this requires an aggregation into binary values, such as *rounded average*, *majority vote*, and *logical OR*. Such techniques, however, lack the precision required by binary data.

To visualize the extent of this problem, we show the impact of straightforward aggregation in Fig. 1(a), which highlights that even the best combination of a local factorization algorithm and an aggregation scheme—here, ASSO (Miettinen et al. 2008) using logical OR—leads to bad reconstructions.

Recently, Dalleiger and Vreeken (2022) showed we can continuously relax BMF into a regularized *binary* matrix factorization problem using linear (rather than Boolean) algebra and proximal gradients, yielding an efficient and scalable approach with state-of-the-art performance. Taking advantage of this relaxation, we propose the FELB algorithm that locally factorizes while centrally, yet privacy-consciously aggregates coefficients using a proximal aggregation, thereby efficiently yielding valid global binary matrices. On our toy example in Fig. 1(b) it achieves a nearly perfect reconstruction.

We show that FELB converges to binary matrices, provide differential privacy guarantees using, e.g., the Gaussian mechanism (Balle and Wang 2018), and we experimentally validate that the utility remains high. Moreover, we demonstrate that FELB outperforms baselines derived via straightforward parallelization of state-of-the-art BMF methods on numerous real-world and synthetic datasets.

In summary, our main contributions are as follows:

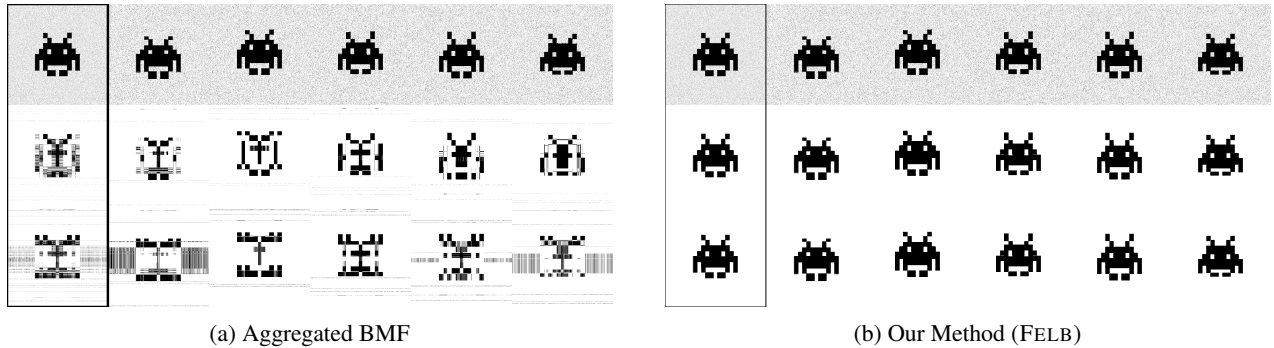


Figure 1: Our method reconstructs data well. Representing 1s as black pixels, for (a) ASSO using logical OR and (b) our novel federated factorization called FELB, we show (top row) the client-data subjected to additive noise, (middle row) the localized reconstructions, and (bottom row) the aggregation-based reconstructions. The left-most column shows the results centralized combination of the data resp. reconstructions of the five clients (columns 2–6).

- We introduce a novel federated proximal-gradient-descent for BMF (FELB).
- We improve over the state-of-the-art in BMF with our adaptive regularization (FELB^{MU}).
- Provide a formal foundation for federated inertial alternating proximal-gradient optimization under non-convex regularization.
- We experimentally show that our methods are both efficient and accurate.

2 Related Work

To the best of our knowledge, there exists no federated BMF algorithms. We therefore primarily discuss the relations to *BMF*, and *federated factorization*, and *federated learning*.

We distinguish two classes of **BMF** methods: First, *discrete optimization-based methods* that use Boolean algebra, such as ASSO (Miettinen et al. 2008) using a set-cover-like approach, GRECOND (Belohlávek and Vychodil 2010), MEBF (Wan et al. 2020) using fast geometric segmentation, or SOFA (Neumann and Miettinen 2020) based on streaming clustering. Second, *continuous optimization-based methods* that use linear algebra for solving the binary matrix factorization problem, introduced by Zhang et al. (2007), and advanced by Araujo, Ribeiro, and Faloutsos (2016) based on thresholding, and by Hess et. al (Hess, Morik, and Piatkowski 2017; Hess and Morik 2017) using a proximal operator. Combining ideas from the two complementary regularization strategies of Hess, Morik, and Piatkowski (2017) and Zhang et al. (2007), Dalleiger and Vreeken (2022) recently removed the need for post-processing via a proximal operator for an elastic-net-based regularizer.

With regards to **federated factorization** in general, ‘parallel’ algorithms for matrix factorization (Yu et al. 2014) as well as binary matrix factorization (Khanna et al. 2013) seek computational efficiency without addressing privacy concerns. Towards matrix factorization for distributed privacy-sensitive data, methods exist for federated matrix factorization (Du et al. 2021) and federated non-negative matrix factorization (Li et al. 2021). These methods, however, are not specialized to Boolean matrices. Here, we close the research gap

by addressing the need for a federated, privacy-preserving binary (or Boolean) matrix factorization algorithm.

Recent advances in **federated learning** involve techniques like FedProx (Li et al. 2020a) and SCAFFOLD (Karimireddy et al. 2020). FedProx, an extension of FedAvg (McMahan et al. 2017), introduces a proximity penalty term to stabilizing the training process across different clients. SCAFFOLD enhances federated learning by correcting client drift using variance reduction techniques, thereby improving convergence rates and model accuracy compared to traditional methods like FedAvg, while ProxSkip (Mishchenko et al. 2022) uses randomization to reduce the computational cost of proximal operators which are significantly more expensive than our operators. Despite these advances, most research focuses on training deep neural networks using stochastic-gradient-based local optimization schemes. These approaches often yield to a slow convergence to suboptimal non-Boolean solutions, if they are deployed to similar non-convex alternating optimization contexts.

3 Federated Proximal Binary Matrix Factorization

Having contextualized our problem, we now formally introduce our federated Boolean matrix factorization scenario, show how we separate our problem into manageable subproblems; describe how to efficiently and solve subproblems in terms of binary matrix factorization relaxation, while preserving privacy; and formally show that we compute a Boolean matrix factorization upon convergence.

The most pronounced difference between traditional and federated Boolean matrix factorization lies in data accessibility. Rather than having all data $A \in \{0, 1\}^{n \times m}$ accessible at one location, the data A is given as (horizontally) partitioned matrices A_1, \dots, A_c over $C \in \mathbb{N}$ clients such that

$$A = [A_1, \dots, A_C]^T,$$

where $A_j \in \{0, 1\}^{n_j \times m}$ and $n = \sum_i n_i$. We aim to discover a single *shared* matrix $\hat{V} \in \{0, 1\}^{k \times m}$ containing shared feature components that are beneficial for all clients. Due to privacy

restrictions, we are however neither permitted to transmit matrices A_i ‘offsite’ (including to any other device), nor are we allowed to be able to draw conclusions about where components belong to. We want to factorize the data $A_i \approx U_i \circ \widehat{V}$ in terms of *local* matrix $U_i \in \{0, 1\}^{\tilde{c} \times k}$ (associating data to components), and one shared *global* matrix $\widehat{V} \in \{0, 1\}^{k \times m}$ (associating features into components). Without the knowledge of U_i , we *cannot* estimate specific attributes of individual users (assuming sufficiently large client datasets). We *can*, however, estimate sets of commonly co-occurring attributes across all clients, e.g. common combinations of genetic markers that are indicative of a disease.

Locally computing U_i for given A_i and \widehat{V} is a regular Boolean matrix factorization. However, computing the *shared* \widehat{V} without access to A_i and U_i is not straightforward. To enable the computing of a shared factor while still preserving privacy, we split the problem into subproblems Φ_i , introducing a local *but shareable* coefficient matrix $V_i \in \{0, 1\}^{k \times m}$. In a nutshell, we estimate a factorization for Φ_i , combine local matrices $V_i \in \{0, 1\}^{k \times m}$ into a shared matrix \widehat{V} , update Φ_i , and repeat. In a nutshell, we seek to optimize

$$\arg \min_{U, V, \widehat{V}} \sum_i \Phi_i(U_i, V_i, \widehat{V}) \quad (1)$$

specifying and solving the subproblems next.

Local Subproblems and Clients

A single subproblem at client $i \in \mathbb{N}$, seeks to optimize $A_i \approx [U_i \circ V_i]_{ab} = \bigvee_{c \in [k]} U_{i,c} V_{i,cb}$, of two low-rank Boolean factor matrices (Miettinen et al. 2008), $U_i \in \{0, 1\}^{n_i \times k}$ (*feature matrix*) and $V_i \in \{0, 1\}^{k \times m}$ (*coefficient matrix*). As this problem is NP-complete (Miettinen et al. 2008), solving it exactly is challenging for each client, even for relatively small matrices.

A major factor contributing to this hardness is the requirement that variables are Boolean. To address this challenge, we essentially relax the Boolean constraint by replacing it with additional penalty terms R and P detailed below. That is, we continuously relax the problem into a *binary matrix factorization* problem

$$\Phi_i(U_i, V_i) \|A_i - U_i V_i\|_F^2 + R(U_i) + R(V_i) + P(V_i) \quad , \quad (2)$$

for relaxed $U_i \in [0, 1]^{n_i \times k}$ and $V_i \in [0, 1]^{k \times m}$ using regular linear algebra. First, To yield the desired Boolean outcomes without constrains, we introduce a *binary-inducing regularizer* $R : \mathbb{R}^{n' \times m'} \rightarrow \mathbb{R}$, enabling efficient gradient-based optimizations. A regularizer that encourages binary solutions combines two elastic-nets (rooted at 0 and 1, resp.) into the almost W-shaped ELB-regularizer

$$R_{\kappa, \lambda}(X) = \sum_{x \in X} \min \{r(x), r(x - \mathbf{1})\} \quad (3)$$

where $r(x) = \kappa \|x\|_1 + \lambda/2 \|x\|_2^2$ (Dalleiger and Vreeken 2022). Second, for faster convergence towards a shared solution, we introduce a proximity penalty $P : \mathbb{R}^{n' \times m'} \rightarrow \mathbb{R}$, which encourages local V_i to remain close to the global model using the distance $P(V_i) = \gamma \|V_i - \widehat{V}\|_F^2$ between them.

Even though now unconstrained, this problem is still challenging due to being non-convex. We solve this joint objective by first splitting it in two subproblems, solving them alternating

$$U_i^{t+1} = \arg \min_U \|A_i - UV_i^{t-1}\|_F^2 + R(U) \quad \text{and}$$

$$V_i^{t+1} = \arg \min_V \|A_i - U_i^{t+1}V\|_F^2 + R(V) + P(V) \quad .$$

Because each individual objective remains a challenge due to the non-convexity, we require an optimization algorithm that is capable of solving such non-convex problems. To this end, we employ the *inertial proximal alternating linear minimization* (iPALM) technique (Pock and Sabach 2016), which will guarantee convergence (Attouch, Bolte, and Svaiter 2013; Bolte, Sabach, and Teboulle 2014) as detailed in Sec. 16.

Proximal Alternating Linear Minimization At the core of iPALM, each regularized objective for U_i and V_i are solved using a proximal gradient approach, which separates loss from regularizer. That is, after taking a gradient step concerning our linear least-squares loss f , e.g., $f(U) \leftarrow \|A_i - UV_i^{t-1}\|_F^2$, we then take a scaled proximal step regarding regularizer to project the gradient towards a feasible Boolean solution and towards a proximity to \widehat{V} for V_i . A proximal operator is the projection

$$\text{prox}_{\eta}^R(X) = \arg \min_Y \frac{1}{2\eta} \|X - Y\|_F^2 + R(X) \quad (4)$$

of the result of the gradient step $x - x\eta \nabla_x f(x)$ for the loss f , into the proximity of a regularized solution $R(X)$. With regards to our regularizer R and P , these proximal problems lend themselves for deriving first-order optimal and efficiently-computable closed-form solutions: The *Boolean proximal operator* for R is element-wise computable

$$\text{prox}_{\eta}^R(X) = \frac{1}{1+\eta\lambda} \text{sign}(X - \theta) \max\{|X - \theta| - \eta\kappa, 0\} \quad , \quad (5)$$

for element-wise indicator $\theta = \mathbb{1}[X_{ij} \leq 1/2]_{ij}$ (Dalleiger and Vreeken 2022), as shown in Apx. A. The *\widehat{V} -proximity proximal operator* for P is simply a weighted average

$$\text{prox}_{\eta}^P(X) = [1 + \eta\gamma]^{-1} (X + \eta\gamma\widehat{V}) \quad . \quad (6)$$

Together, they yield the alternating update rules

$$U_i^{t+1} = \text{prox}_{\eta U_i}^R (U_i^t - \eta_{U_i}^t \nabla_{U_i}^t \|A_i - U_i^t V_i^t\|_F^2)$$

$$V_i^{t+1} = \text{prox}_{\eta V_i}^P \text{prox}_{\eta V_i}^R (V_i^t - \eta_{V_i}^t \nabla_{V_i}^t \|A_i - U_i^{t+1} V_i^t\|_F^2) \quad . \quad (7)$$

To apply these rules, we require step sizes, utilizing linear nature of the loss, we propose two alternatives: first we use the gradient Lipschitz constant L for $\eta = 1/L$, yielding the update rule for FELB. Second we employ Lee and Seung (2000)’s *multiplicative update rule* (MU) for NMF with step size matrices $\eta_{U_i}^t = U_i \oslash U_i V_i V_i^T$ and $\eta_{V_i}^t = V_i \oslash U_i^T U_i V_i$ using the Hadamard division \oslash , containing individual step sizes for all elements in U_i and V_i , yielding FELB^{MU}.

Global Objective and Server

Now having established our per client subproblems, we now combine the local subobjectives into one global objective $\Phi(U, V, \widehat{V}) = \sum_i \Phi_i(U_i, V_i, \widehat{V})$ which is

$$\sum_i \|A_i - U_i V_i\|_{\mathbb{F}}^2 + R(U_i) + R(V_i) + R(\widehat{V}) + P(V_i), \quad (8)$$

focusing on shared coefficients components \widehat{V} . To estimate the shared matrix \widehat{V} independent of all data matrices A_i and local basis matrices U_i , we have to combine V_i matrices. In federated learning, this is often done by aggregating all V_i as the average \widehat{V} . However, doing so here does not necessarily yield valid results: naïve averaging results in aggregates that are far from being binary, thus hindering or even preventing convergence. Addressing this aggregation problem, we aim to result in a Boolean matrix, for which we iteratively project the aggregate towards a valid Boolean values

$$\widehat{V} \leftarrow \arg \min_{\widehat{V}} \sum_i 1/2 \|\widehat{V} - V_i\|_{\mathbb{F}}^2 + R(\widehat{V}), \quad (9)$$

for which we employ a proximal aggregation yielding the update-step $\widehat{V} \leftarrow \text{prox}_{\frac{R}{c}}^{\frac{1}{c}} \sum V_i$. To theoretically guarantee *differential privacy*, clients may further distort the matrices V_i before transmission, as described next.

Guaranteeing Differential Privacy

The proposed aggregation approach only shares coefficient matrices, so that no direct relationships between observations are shared. An attacker or a curious server can, however, attempt to infer private data from coefficients V_i . Aiming to prevent this, we guarantee differential privacy using an additive noise mechanisms, where, in a nutshell, each client adds noise before it transmits V_i to the server. We consider the Bernoulli, Gaussian, and Laplacian mechanisms, which only differ in the noise distribution. Using a Gaussian mechanism, we achieve (ϵ, δ) -differential privacy as follows.

Definition 3.1 (Dwork, Roth et al. (2014)). For $\epsilon, \delta > 0$, a randomized algorithm $\mathcal{A} : \mathcal{X} \rightarrow \mathcal{Y}$ is (ϵ, δ) -differentially private (DP) if

$$P(\mathcal{A}(X) \in S) \leq e^\epsilon P(\mathcal{A}(X') \in S) + \delta$$

holds for each subset $S \subset \mathcal{Y}$ and for all pairs of neighboring inputs X, X' .

Applying Gaussian noise with 0 mean and σ variance to the local coefficients V_i before sending ensures (ϵ, δ) -DP (Balle and Wang 2018) for $\sigma = \Delta \epsilon^{-1} \sqrt{2 \log(5/(4\delta))}$, where $\Delta = \sup_{X, X'} \|\mathcal{A}(X) - \mathcal{A}(X')\|$ is the sensitivity of \mathcal{A} . To ensure bounded sensitivity, we clip all V_i with clipping threshold $\theta > 1$ (Noble, Bellet, and Dieuleveut 2022). Similarly, adding 0-mean $\Delta \epsilon^{-1}$ -variance Laplacian noise achieves $(\epsilon, 0)$ -DP (Dwork et al. 2006).

Convergence Analysis

Having ensured differential privacy, we summarize our algorithm. We call the combination of this proximal aggregation with local proximal-gradient optimization steps the FELB algorithm, detailed in Alg. 1: Local factors U_i, V_i are initialized

Algorithm 1: Federated Binary Matrix Factorization with FELB

Input: distributed target matrices A^1, \dots, A^C , component-count k

Output: local feature matrices U_1, \dots, U_C , global coefficient matrix V

```

1 initialize  $U_i, V_i$  for  $i \in [C]$  uniformly at random
2 Locally at client  $i$  in iteration  $t$  do
3    $U_i \leftarrow \text{prox}_{\eta U_i}^{R_i}(U_i - \eta_{U_i} \nabla_{U_i} \|A_i - U_i V_i\|_{\mathbb{F}}^2)$ 
4    $V_i \leftarrow \text{prox}_{\eta V_i}^{R_i}(V_i - \eta_{V_i} \nabla_{V_i} \|A_i - U_i V_i\|_{\mathbb{F}}^2)$ 
5    $V_i \leftarrow \text{prox}_{\eta V_i}^P(V_i)$ 
6   if  $t \bmod b = 0$  then
7     if is differentially private then
8        $V_i \leftarrow V_i \oplus N, N_{ab} \sim \mathcal{N}(0, \sigma)$ 
9     transmit  $V_i$  to the server
10    receive  $\widehat{V}$  from the server
11    let  $V_i \leftarrow \widehat{V}$ 
12 At server do
13   receive  $V_1, \dots, V_C$ 
14   aggregate  $\widehat{V} \leftarrow \text{prox}_{\frac{R}{c}}^{\frac{1}{c}} \left( \frac{1}{c} \sum_{i=1}^C V_i \right)$ 
15   transmit  $\widehat{V}$  to each client
16 return  $U, \widehat{V}$ 

```

uniformly at random (line 1), and at each client in round t (line 2), we update the local factor matrices (lines 4 and 5). Every b rounds, we transmit the local matrices V_i to the server (line 7). At this point, each client may choose to preserve differential privacy. The server, receives all local coefficients V_i (line 11), averages the matrices, and applies the proximal-operator (line 12). The aggregate is then transmitted to all clients (line 13). Upon receiving the aggregate (lines 8 and 9), each client continues with the next optimization round.

Next, to formally ascertain that Alg. 1 solves our problem, we show that the algorithm converges with Thm. 3.2, and achieves Boolean coefficients in the limit with Thm. 3.3.

Theorem 3.2 (Convergence). *For the sequence generated by Alg. 1 $\{z^t \triangleq (\{U_i^t\}_i, \{V_i^t\}_i, \widehat{V}^t)\}_{t \in \mathbb{N}}$, the objective function $\Phi(z^t)$ converges to a stable solution $\Phi(z^t) \rightarrow \widehat{\Phi}$ if $t \rightarrow \infty$.*

Proof. (Sketch, full proof in Apx. B). We show the objective's convergence to a stable solution Φ^* by initially establishing the convergence of each client, where we observe a *sufficient reduction* in local objectives, as well as a *bounded dissimilarity* to \widehat{V} . Leveraging this, we establish global convergence by showing that the global loss gradient is bounded by a *diminishing term*, showing that $\Phi(z^t)$ approaching a constant $\widehat{\Phi}$ as t tends to infinity. \square

Theorem 3.3 (Boolean Convergence). *If λ^t is a monotonically increasing sequence with $\lambda^{t-1} \leq \lambda^t$, $\lim \lambda^t \rightarrow \infty$, and $\lambda^t - \lambda^{t-1} \leq \infty$, then V_1^T, \dots, V_c^T and \widehat{V}^T from the sequence generated by Alg. 1 converges as $\lim_{T \rightarrow \infty} \text{dist}(\widehat{V}^T, \{0, 1\}) \rightarrow 0$ to a Boolean matrix.*

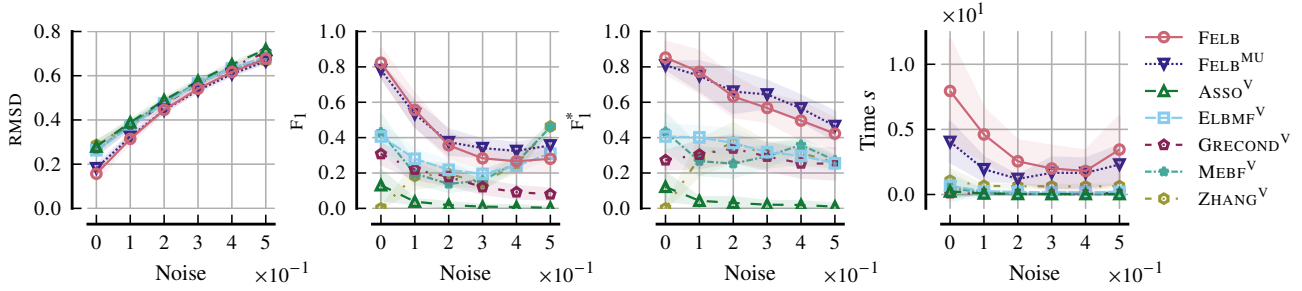


Figure 2: FELB and FELB^{MU} are robust against noise. We show the loss, recall, similarity, and elapsed runtime (s/C) for synthetic data with varying levels of destructive XOR noise.

Proof. (Sketch, full proof in Apx. B). Since gradients are bounded and diminish, we only need to show that the proximal operator returns Boolean solutions in the limit. As our gradients are Lipschitz continuous, bounded, and ensured to converge to a stable solution, our scaled proximal operator projects values onto Boolean results, for a monotonically increasing regularizer rate λ' that approaches infinity in the limit, guaranteeing a stable Boolean convergence regardless of communication rounds. \square

Turning traditional BMF into FedBMF

Given that there exist no federated matrix factorization algorithms tailored to binary data, we compare our approaches to local BMF methods, whose outcomes are then partially transmitted to a central location and collectively aggregated, following established ad-hoc federation strategies (Kamp 2019). In particular, we adapt the localized algorithms, covering the state of the art in the method families (1) *cover-based Boolean matrix factorizations* (ASSO, Miettinen et al. (2008); GRECOND, Belohlávek and Vychodil (2010); MEBF, Wan et al. (2020)) and (2) *relaxation-based binary matrix factorizations* (ZHANG, Zhang et al. (2007); and ELBMF, Dalleiger and Vreeken (2022)), to factorize *distributed* matrices—factorizing locally and aggregating the coefficient matrices centrally, replacing the local coefficients. Leveraging the following aggregations, we summarize the BMF federation scheme in Apx. C Alg. 2. To ensure binary results, we employ three *aggregation strategies* that maintain valid matrices

$$\text{Rounded Average} \quad \lfloor C^{-1} \sum_{c \in [C]} V^c \rfloor \quad (10)$$

$$\text{Majority Vote} \quad \left[\sum_{c \in [C]} V_{ij}^c \geq C/2 \right]_{ij} \quad (11)$$

$$\text{Logical OR} \quad \bigvee V^1, \dots, V^C \quad (12)$$

We now describe our diverse set of experimental setups. First, we ascertain that FELB works reliably on synthetic data. Second, we empirically assess the differential-privacy properties of FELB. And third, we verify that FELB performs well on diverse real-world datasets drawn from four different scientific areas. To quantify the results, we report the *root mean squared deviation* (RMSD) and the F_1 score between data and reconstruction, as well as the runtime in seconds.

4 Experiments

We implement FELB in the Julia language and run experiments on 32 CPU Cores of an AMD EPYC 7702 or one NVIDIA A40 GPU, reporting wall-clock time in seconds. We provide the source code, datasets, synthetic dataset generator,¹ and additional information regarding reproducibility in Apx. E. In all experiments, we limit each algorithm run to 12h in total. We quantify the performance of *federated* ASSO, GRECOND, ELBMF, MEBF, FELB, and FELB^{MU} in terms of loss, recall, similarity, and runtime, reporting results for *majority voting* in the following, as it has superior performance to *rounded averaging* and *logical*, as shown in Apx. F.

Experiments on Synthetic Data

In these experiments, we answer the following questions:

- Q1 How robust are the algorithms wrt. noise?
- Q2 How scalable are they with increasing client counts?
- Q3 How well do they perform under differential privacy?

To answer these, we need a controlled test environment. We construct this by sampling random binomial-noise matrices, into which we insert randomly generated, densely populated ‘tiles’ containing approximately 90% with 1s. To highlight trends, rather than random fluctuations, we report the mean and confidence intervals of 10 randomized trials.

Robustness regarding Noise To study the impact of noise on the quality of reconstructions, we generated synthetic matrices with varying degrees of destructive XOR noise, ranging from 0% (no noise, consisting solely of high-density tiles) to a maximum of 50% (completely random noise). Employing a fixed number of 10 clients, we applied federated ASSO, GRECOND, MEBF, ELBMF, and ZHANG, alongside FELB and FELB^{MU} to each dataset.

We present RMSD, F_1 score (re signal and noise data), F_1^* score (re signal), and runtime in Fig. 2: We see that reconstruction quality declines with increasing noise, yet FELB and FELB^{MU} achieve the best reconstructions across the board even at high noise levels. We see that RMSD and F_1 follow a similar trend across all methods, yet our methods consistently outperform the rest. However, if we regard only the interesting data signal with F_1^* , we see that FELB and FELB^{MU}

¹<https://doi.org/10.5281/zenodo.14501661>

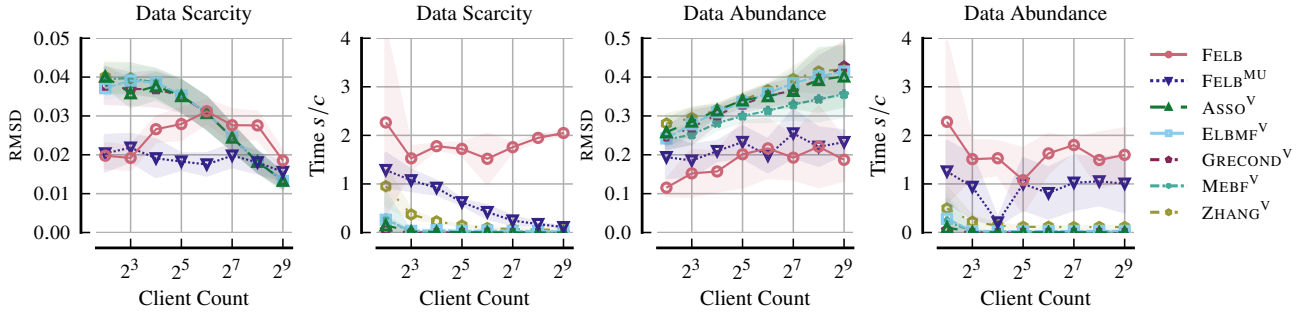


Figure 3: FELB and FELB^{MU} perform well across various client counts, showing RMSD and runtime (s/C). For *data scarcity*, we fix the data size and an increase number of clients. For *data abundance* we grow data while increasing the number of clients.

are the only algorithms that result in good reconstructions of the ground-truth signal, even under high noise. This shows the ability of FELB and FELB^{MU} to discern pure noise from meaningful signal. While the runtime of ASSO, GRECOND, MEBF, ZHANG, and ELBMF is slightly faster in Fig. 2 (right), FELB^{MU}'s and FELB's runtime reduces with increasing noise.

Scalability regarding Clients Next, we analyze the scalability of federated ASSO, GRECOND, ELBMF, MEBF, and ZHANG under *majority voting*, as well as of FELB and FELB^{MU}, for varying numbers of clients, considering two contrasting scenarios of *scarce* and *abundant* data. In both cases, we generate and uniformly distribute synthetic data to a number of clients, depicting results in Fig. 3.

To create *data scarcity*, we fix the dataset size to 2^{16} and increase the number of clients from 2^2 to 2^9 , thus iteratively reducing the sample count per client. In Fig. 3 (left), we observe that our methods scale well to low-sample scenarios and deliver the best performance. The MU update rule outperforms the competitors. The runtime of post-hoc federated methods ASSO, GRECOND, MEBF, ZHANG, and ELBMF is lower since they only perform a single optimization epoch. These methods slightly outperform FELB and FELB^{MU} only in tiny data scenarios where the estimator-variance is high, while the FELB^{MU} significantly outperforms all methods and is notably faster than FELB.

To evaluate under *data abundance*, we scale the number of samples by increasing the number of clients from 2^2 to 2^9 , maintaining a constant sample count of 500 per client. In Fig. 3 (right), we observe that our methods scale well with an increased number of clients. With more data, FELB using Lipschitz steps slightly outperforms the MU steps in RMSD, and both methods exhibit comparable runtime trends. The runtime of post-hoc federated methods ASSO, GRECOND, MEBF, ZHANG, and ELBMF remains lower, as they compute only one local optimization epoch.

Performance under Privacy To empirically ascertain the effect of differential-privacy on the loss, we add noise to the transmitted factor matrices according to various noise mechanisms. Specifically, we study the effect of additive clipped or regular Laplacian and Gaussian, as well as xor Bernoulli noise mechanisms, as depicted in Fig. 4 and Apx. F, for vary-

ing $0 \leq \epsilon \leq 2$ and fixed $\delta = 0.05$. Because ASSO, MEBF, GRECOND, ZHANG, and ELBMF return Boolean matrices, we subject these only to xor noise, rather than additive noise, to retain Boolean matrices. The results in Fig. 4 show that both FELB and FELB^{MU} exhibit similar performance across various noise models, while FELB^{MU} is most robust. The plots display three phases: In the low- ϵ domain, there is almost no performance deterioration, followed by a steep, hockey-stick-like descent which eventually stabilizes in the high- ϵ range. We note an increasing 'sharpness' of the hockey-stick-phase under clipping, showing less smooth reactions to privacy adjustments for both mechanisms.

Experiments on Real-World Data

Having established the efficiency and precision of our method on synthetic data, we proceed to assess its effectiveness on real-world datasets. For this, we curated a diverse set of 8 real-world datasets spanning four distinct domains. To explore **recommendation systems**, we include *Goodreads* (Kotkov et al. 2022) for books and *Movielens* (Harper and Konstan 2015) and *Netflix* (Netflix, Inc. 2009) for movies, where user ratings ≥ 3.5 are binarized to 1. In **life sciences**, we use *TCGA* (Institute 2005) for cancer genomics, *HPA* (Bakken et al. 2021; Sjöstedt, Zhong, and et. al 2020) for single-cell proteomics, and *Genomics* (Oleksyk, Gonçalo, and et. al 2015) for mutation data. *TCGA* marks gene expressions in the top 95% quantile as 1, while *HPA* designates observed RNA in cells as 1. For **social science**, we analyze poverty (*Pov*) and income (*Inc*) using the ACS (U.S. Census Bureau 2023) dataset, binarizing with one-hot encoding utilizing Folktables (Ding et al. 2021). In **natural language processing**, we study higher-order word co-occurrences in ArXiv cs.LG abstracts (Collaboration 2023). Each paper abstract is a row with columns marked 1 if the corresponding word is in the vocabulary, containing lemmatized, stop-word-free words with a minimum frequency of 1 ‰. We summarize dataset extents, density, and chosen component counts in Apx. E, Tbl. 2. Since the number of clients (e.g., Hospitals) is expected to be small, we limit the federation to a reasonable $C = 50$ clients, on which we compare federated methods ASSO, GRECOND, MEBF, ELBMF, and ZHANG, as well as FELB, and FELB^{MU} across all real-world datasets,

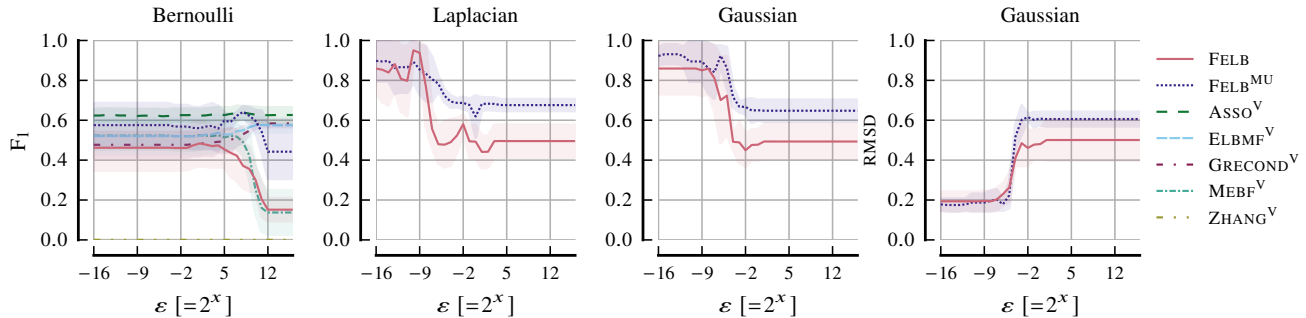


Figure 4: FELB and FELB^{MU} achieve accurate yet differentially private reconstructions. For synthetic data, we subject algorithms to different noise mechanisms: Bernoulli, Laplacian, and Gaussian noise.

Table 1: FELB and FELB^{MU} consistently perform well. We illustrate the F₁ of ASSO, GRECOND, MEBF, ELBMF, and ZHANG under voting aggregation, as well as federated FELB, and FELB^{MU} on 8 real-world data across 50 clients. We highlight the best scoring algorithm with **bold**, the second best with underline, and timeouts by a dash –.

Dataset	ASSO ^V	MEBF ^V	GRECOND ^V	ZHANG ^V	ELBMF ^V	FELB ^{MU}	FELB
ACS Inc	0.388	0.108	0.690	0.000	0.000	<u>0.585</u>	0.328
ACS Pov	<u>0.692</u>	–	0.797	0.000	0.217	<u>0.638</u>	0.517
cs.LG	–	0.000	0.068	0.000	0.000	<u>0.057</u>	0.006
Goodreads	–	0.000	0.017	–	0.000	0.125	<u>0.059</u>
HPA Brain	–	0.642	–	0.000	0.000	<u>0.007</u>	0.000
Movielens	–	0.017	–	–	0.000	0.193	<u>0.163</u>
Netflix	–	0.010	–	–	0.000	0.197	<u>0.144</u>
TCGA	0.039	0.055	0.007	0.000	0.000	0.414	<u>0.402</u>
<i>Avg. Rank</i>	4.750	3.75	3.375	5.125	4.500	1.625	<u>2.750</u>

synchronizing after every $b = 10$ local optimization rounds.

In Tbl. 1, we present the F₁ between the reconstruction and the data matrix, where – indicate missing data due to time limits. Our results show that FELB and FELB^{MU} exhibit best-in-class performance, consistently ranking as the best or second-best algorithms. This performance gap is evident across all datasets except for the *HPA* dataset, where MEBF, a method designed with similar data types in mind, outperforms the others, and the *ACS Pov* dataset, where GRECOND leads. Notably, since clients of ELBMF and ZHANG diverge significantly, they often aggregate into a no-consensus 0-only global model matrix, thus showing low accuracy. Although they perform only a single optimization round per client, we see that ASSO, GRECOND, and MEBF do not finish on medium to large datasets. Additionally, we show the RMSD in Apx. F, where FELB and FELB^{MU} are on top, and compare client-server communication frequencies in Apx. F, demonstrating the strength of FELB and resp. FELB^{MU}.

5 Discussion and Conclusion

We introduced the federated proximal-gradient-based FELB for BMF tasks, showed its convergence to a binary outcome in theory, and demonstrated its efficacy in experimental practice. We provided a variant called FELB^{MU}, whose practical performance outcompetes FELB on many real-world datasets,

especially under rare synchronizations. Although FELB and FELB^{MU} perform consistently well, both are first-of-their-kind federated BMF algorithms. As such, they leave ample room for further research.

Limitations Our research focuses on learning from private Boolean data generated by similar sources at a few research centers, thus we concentrate on suitable experiments and abstain from distant but related problems, such as learning with millions of heterogeneous clients. Further, we experimentally demonstrate the practical limitations of our methods extensively, extending this discussion in Apx. G.

Future Work includes extending FELB to allow for heterogeneous clients and data distributions, adapting our methods to learn from varied data distributions and characteristics. Additionally, we plan to explore large-scale federations, drawing inspiration from frameworks like Scaffold (Karimireddy et al. 2020) and FedProx (Li et al. 2020b) for efficient client sampling, variance controlling, and formal limits to client dropout resilience. Furthermore, we intend to investigate personalized federated learning techniques to improve the reconstructions in case of varied data sources. Finally, we plan to move beyond Boolean data and seek explore the potential of allowing partial sharing of a subset of the client components V_i to allow for multi-source multi-modal federated learning to improve model performance and generality.

Acknowledgements

This research is supported by the Wallenberg AI, Autonomous Systems and Software Program (WASP) funded by the Knut and Alice Wallenberg Foundation. Michael Kamp received support from the Cancer Research Center Cologne Essen (CCCE).

References

- Araujo, M.; Ribeiro, P. M. P.; and Faloutsos, C. 2016. Fast-Step: Scalable Boolean Matrix Decomposition. In *Advances in Knowledge Discovery and Data Mining - 20th Pacific-Asia Conference, PAKDD 2016, Auckland, New Zealand, April 19-22, 2016, Proceedings, Part I*, volume 9651 of *Lecture Notes in Computer Science*, 461–473. Springer.
- Attouch, H.; Bolte, J.; and Svaiter, B. F. 2013. Convergence of descent methods for semi-algebraic and tame problems: proximal algorithms, forward-backward splitting, and regularized Gauss-Seidel methods. *Math. Program.*, 137(1-2): 91–129.
- Bakken, T. E.; Jorstad, N. L.; Hu, Q.; and et. al, B. B. L. 2021. Comparative cellular analysis of motor cortex in human, marmoset and mouse. *Nature*, 598: 111 – 119.
- Balle, B.; and Wang, Y.-X. 2018. Improving the gaussian mechanism for differential privacy: Analytical calibration and optimal denoising. In *International Conference on Machine Learning*, 394–403. PMLR.
- Belohlávek, R.; and Vychodil, V. 2010. Discovery of optimal factors in binary data via a novel method of matrix decomposition. *Journal of Computer and System Sciences*, 76(1): 3–20.
- Bolte, J.; Sabach, S.; and Teboulle, M. 2014. Proximal alternating linearized minimization for nonconvex and nonsmooth problems. *Mathematical Programming*, 146(1-2): 459–494.
- Collaboration, A. 2023. arXiv dataset and metadata of 1.7M+ scholarly papers across STEM.
- Dalleiger, S.; and Vreeken, J. 2022. Efficiently Factorizing Boolean Matrices using Proximal Gradient Descent. *Advances in Neural Information Processing Systems*.
- Ding, F.; Hardt, M.; Miller, J.; and Schmidt, L. 2021. Retiring Adult: New Datasets for Fair Machine Learning. *Advances in Neural Information Processing Systems*, 34.
- Du, Y.; Zhou, D.; Xie, Y.; Shi, J.; and Gong, M. 2021. Federated matrix factorization for privacy-preserving recommender systems. *Applied Soft Computing*, 111: 107700.
- Dwork, C.; McSherry, F.; Nissim, K.; and Smith, A. 2006. Calibrating noise to sensitivity in private data analysis. In *Theory of Cryptography: Third Theory of Cryptography Conference, TCC 2006, New York, NY, USA, March 4-7, 2006. Proceedings 3*, 265–284. Springer.
- Dwork, C.; Roth, A.; et al. 2014. The algorithmic foundations of differential privacy. *Foundations and Trends® in Theoretical Computer Science*, 9(3–4): 211–407.
- Golub, G. H.; and Loan, C. F. V. 1996. *Matrix Computations, Third Edition*. Johns Hopkins University Press. ISBN 978-0-8018-5414-9.
- Haddad, A.; Shamsi, F.; Zhu, L.; and Najafizadeh, L. 2018. Identifying Dynamics of Brain Function Via Boolean Matrix Factorization. In *52nd Asilomar Conference on Signals, Systems, and Computers, ACSSC 2018, Pacific Grove, CA, USA, October 28-31, 2018*, 661–665. IEEE.
- Harper, F. M.; and Konstan, J. A. 2015. The MovieLens Datasets: History and Context. *ACM Trans. Interact. Intell. Syst.*, 5(4).
- Hess, S.; and Morik, K. 2017. C-SALT: Mining Class-Specific ALTerations in Boolean Matrix Factorization. In *Machine Learning and Knowledge Discovery in Databases - European Conference, ECML PKDD 2017, Skopje, Macedonia, September 18-22, 2017, Proceedings, Part I*, volume 10534 of *Lecture Notes in Computer Science*, 547–563. Springer.
- Hess, S.; Morik, K.; and Piatkowski, N. 2017. The PRIMPING routine - Tiling through proximal alternating linearized minimization. *Data Mining and Knowledge Discovery*, 31(4): 1090–1131.
- Ignatov, D. I.; Nenova, E.; Konstantinova, N.; and Konstantinov, A. V. 2014. Boolean Matrix Factorisation for Collaborative Filtering: An FCA-Based Approach. In *Artificial Intelligence: Methodology, Systems, and Applications - 16th International Conference, AIMSA 2014, Varna, Bulgaria, September 11-13, 2014. Proceedings*, volume 8722 of *Lecture Notes in Computer Science*, 47–58. Springer.
- Institute, N. C. 2005. The Cancer Genome Atlas Program (TCGA).
- Kamp, M. 2019. *Black-box parallelization for machine learning*. Ph.D. thesis, Universitäts-und Landesbibliothek Bonn.
- Karimireddy, S. P.; Kale, S.; Mohri, M.; Reddi, S.; Stich, S.; and Suresh, A. T. 2020. Scaffold: Stochastic controlled averaging for federated learning. In *International Conference on Machine Learning*, 5132–5143. PMLR.
- Khanna, R.; Zhang, L.; Agarwal, D.; and Chen, B.-C. 2013. Parallel matrix factorization for binary response. In *2013 IEEE International Conference on Big Data*, 430–438. IEEE.
- Kotkov, D.; Medlar, A.; Maslov, A.; Satyal, U. R.; Neovius, M.; and Glowacka, D. 2022. The Tag Genome Dataset for Books. In *Proceedings of the 2022 ACM SIGIR Conference on Human Information Interaction and Retrieval (CHIIR '22)*. ISBN 978-1-4503-9186-3/22/03.
- Lee, D. D.; and Seung, H. S. 1999. Learning the parts of objects by non-negative matrix factorization. *Nature*, 401(6755): 788–791.
- Lee, D. D.; and Seung, H. S. 2000. Algorithms for Non-negative Matrix Factorization. In *Advances in Neural Information Processing Systems 13, Papers from Neural Information Processing Systems (NIPS) 2000, Denver, CO, USA*, 556–562. MIT Press.
- Li, T.; Sahu, A. K.; Zaheer, M.; Sanjabi, M.; Talwalkar, A.; and Smith, V. 2020a. Federated optimization in heterogeneous networks. *Proceedings of Machine Learning and Systems*, 2: 429–450.
- Li, T.; Sahu, A. K.; Zaheer, M.; Sanjabi, M.; Talwalkar, A.; and Smith, V. 2020b. Federated Optimization in Heterogeneous Networks. In Dhillon, I. S.; Papailiopoulos, D. S.;

- and Sze, V., eds., *Proceedings of Machine Learning and Systems 2020, MLSys 2020, Austin, TX, USA, March 2-4, 2020*. mlsys.org.
- Li, Z.; Ding, B.; Zhang, C.; Li, N.; and Zhou, J. 2021. Federated matrix factorization with privacy guarantee. *Proceedings of the VLDB Endowment*, 15(4).
- Liang, L.; Zhu, K.; and Lu, S. 2020. BEM: Mining Coregulation Patterns in Transcriptomics via Boolean Matrix Factorization. *Bioinformatics*, 36(13): 4030–4037.
- McMahan, B.; Moore, E.; Ramage, D.; Hampson, S.; and y Arcas, B. A. 2017. Communication-efficient learning of deep networks from decentralized data. In *International Conference on Artificial Intelligence and Statistics*, 1273–1282. PMLR.
- Miettinen, P.; Mielikäinen, T.; Gionis, A.; Das, G.; and Manilla, H. 2008. The Discrete Basis Problem. *IEEE Transactions on Knowledge and Data Engineering*, 20(10): 1348–1362.
- Mishchenko, K.; Malinovsky, G.; Stich, S.; and Richtarik, P. 2022. ProxSkip: Yes! Local Gradient Steps Provably Lead to Communication Acceleration! Finally! In Chaudhuri, K.; Jegelka, S.; Song, L.; Szepesvari, C.; Niu, G.; and Sabato, S., eds., *Proceedings of the 39th International Conference on Machine Learning*, volume 162 of *Proceedings of Machine Learning Research*, 15750–15769. PMLR.
- Nesterov, Y. 2004. *Introductory Lectures on Convex Optimization*. Springer US. ISBN 9781441988539.
- Netflix, Inc. 2009. Netflix Prize.
- Neumann, S.; and Miettinen, P. 2020. Biclustering and Boolean Matrix Factorization in Data Streams. *Proceedings of the VLDB Endowment*, 13(10): 1709–1722.
- Noble, M.; Bellet, A.; and Dieuleveut, A. 2022. Differentially private federated learning on heterogeneous data. In *International Conference on Artificial Intelligence and Statistics*, 10110–10145. PMLR.
- Oleksyk, T. K.; Gonçalo, A.; and et. al, R. D. 2015. A global reference for human genetic variation. *Nature*, 526: 68–74.
- Paatero, P.; and Tapper, U. 1994. Positive matrix factorization: A non-negative factor model with optimal utilization of error estimates of data values. *Environmetrics*, 5(2): 111–126.
- Parikh, N.; and Boyd, S. P. 2014. Proximal Algorithms. *Foundations and Trends in Optimization*, 1(3): 127–239.
- Pock, T.; and Sabach, S. 2016. Inertial Proximal Alternating Linearized Minimization (iPALM) for Nonconvex and Nonsmooth Problems. *SIAM Journal on Imaging Sciences*, 9(4): 1756–1787.
- Sjöstedt, E.; Zhong, W.; and et. al, L. F. 2020. An atlas of the protein-coding genes in the human, pig, and mouse brain. *Science*, 367(6482): eaay5947.
- U.S. Census Bureau. 2023. Census. U.S. Department of Commerce.
- Wan, C.; Chang, W.; Zhao, T.; Li, M.; Cao, S.; and Zhang, C. 2020. Fast and Efficient Boolean Matrix Factorization by Geometric Segmentation. In *The Thirty-Fourth AAAI Conference on Artificial Intelligence, AAAI 2020, New York, NY, USA, February 7-12, 2020*, 6086–6093. AAAI Press.
- Yu, H.-F.; Hsieh, C.-J.; Si, S.; and Dhillon, I. S. 2014. Parallel matrix factorization for recommender systems. *Knowledge and Information Systems*, 41: 793–819.
- Zhang, Z.; Li, T.; Ding, C. H. Q.; and Zhang, X. 2007. Binary Matrix Factorization with Applications. In *Proceedings of the 7th IEEE International Conference on Data Mining (ICDM 2007), October 28-31, 2007, Omaha, Nebraska, USA*, 391–400. IEEE Computer Society.

Supplementary Material

In this Appendix, we provide supplementary information

- regarding the convergece in Apx. B,
- regarding the federation of baseline BMF methods in Apx. C,
- regarding dataset used in our experiments Apx. D,
- regarding reproducibility of our experiments in Apx. E,
- regarding limitations in Apx. G.

Furthermore, we provide additional experimental results regarding

- post-hoc aggregations in Apx. F,
- empirical convergence in Apx. F,
- client drift in Apx. F
- differential privacy in Apx. F, and
- additional real-world performance evaluations in Apx. F.

A Binary Elastic-Net Proximal Operator

For the ambiguity-free minimum

$$\min_{\leq} \{A, B\} = \begin{cases} A & A \leq B \\ B & A > B \end{cases},$$

we define the binary elastic-net as

$$R(X) = \min_{\leq} \{r(X), r(X-1)\} \quad \text{where} \quad r(X) = \lambda/2 \|X\|_2^2 + \kappa \|X\|_1.$$

To obtain a proximal operator we need the minimizer of

$$\arg \min_X R(X) + \frac{1}{\eta} \|X - Y\|_F^2.$$

Rather than scaling the proximity, we equivalently (inversely) scale the regularization, yielding

$$\text{prox}_{\eta}^R(X) = \arg \min_X R(X) + 1/\eta \|X - Y\|_F^2.$$

While our iPALM-based approach relies on a single step size η per gradient matrix (i.e., η_U or η_V), our MUR-based approach uses one step size *per cell* of the gradient matrix, thereby adaptively adjusting the optimization rate. To cover both cases, we Hadamard multiply both terms $R(X) + 1/\eta \|X - Y\|_F^2$, by η , which corresponds to an adaptive scaling of our parameters

$$\lambda_{ij} \triangleq \lambda \eta_{ij} \quad \text{and} \quad \kappa_{ij} \triangleq \kappa \eta_{ij}.$$

Noticing that this equation is piecewise-symmetric and convex, we separately solve each case in turn. For $R(X) \leq R(X-1)$, we obtain the regular elastic-net operator (Parikh and Boyd 2014)

$$\text{prox}^R(X) = [1 + \lambda]^{-1} [\text{sign}(X) \max(|X| \ominus \kappa, 0)]$$

Being element-wise solvable, we notice that this case is taken whenever $X_{ij} \leq 1/2$. Similarly for $R(X) > R(X-1)$, we obtain

$$\text{prox}^R(X) = [1 + \lambda]^{-1} [\text{sign}(X-1) \max(|X-1| \ominus \kappa, 0)].$$

We combine both into the proximal operator

$$\text{prox}^R(X) = [1 + \lambda]^{-1} [\text{sign}(X - \theta_X) \max(|X - \theta_X| \ominus \kappa, 0)]_{ij}$$

where $\theta_X = (2 \text{sign}(X \ominus 1/2) \ominus 1)$.

Now considering the case when $(R(X) = R(X-1))$. Although we ruled-out this ambiguity in practice, in theory, we would otherwise need to consider the convex hull of the differentials in both cases, as (informally)

$$\partial \text{prox}^{R,=} (X) = \text{conv} \{ \partial \text{prox}^{R; < 1/2} (X), \partial \text{prox}^{R; > 1/2} (X) \}$$

which we do in our general convergence analysis. However, by preventing ambiguity, we will simply project towards the more prevalent value of 0 in P in case of a rare tie in practice.

B Convergence

Here, we establish the convergence properties of Algorithm 1. We begin by showing in Theorem B.1 that the objective function of the algorithm converges to a stable solution in the limit. To this end, we leverage the local convergence of each client, as proven with Lem. B.3 (Apx. B), to demonstrate a sufficient reduction in the *global* objective function values. By combining these results, we establish the global convergence of the objective function. Building upon this, we moreover prove in Theorem B.2 that the algorithm converges to Boolean matrices. We establish conditions under which the sequences of matrices converge to binary solutions, demonstrating that both the gradient and proximal operator converge to binary solutions, thereby ensuring the stability of Boolean solutions at both the global and local levels. The outline of our proof is as follows.

1. We show the convergence of Alg. 1 in Thm. B.1.
2. We show that Alg. 1 converges to Boolean matrices with Thm. B.2.
3. We show the convergence of each client in Alg. 1 to a stable solution with Lem. B.3.

Theorem B.1 (Convergence of Alg. 1 (restated)). *For the sequence generated by Alg. 1 $\{z^t \triangleq (\{U_i^t\}_i, \{V_i^t\}_i, \bar{V}^t)\}_{k \in \mathbb{N}}$, the objective function $\Phi(z^t)$ converges to a stable solution $\Phi(z^t) \rightarrow \hat{\Phi}$ if $t \rightarrow \infty$.*

Proof. To show that the objective convergence to a stable solution $\Phi(z^t) \rightarrow \Phi^*$ when $t \rightarrow \infty$, we first show that each client convergence in Lem. B.3, where we observe a sufficient reduction in $\Phi_i(z_i^{t+1}) \leq \Phi_i(z_i^t) - \rho_i \|z_i^{t+1} - z_i^t\|_F^2$ for some constant ρ_i . Using this property we can show the global convergence as follows.

$$\begin{aligned} \Phi(z^{t+1}) &= \sum_i \Phi_i(z_i^{t+1}) \\ &\leq \sum_i \Phi_i(z_i^t) - \rho_i \|\nabla_i \Phi_{z_i^t}(z_i^t)\|_F^2 \\ &\leq \Phi(z^t) - \sum_i \rho_i \|z_i^{t+1} - z_i^t\|_F^2 \\ &\leq \Phi(z^t) - \rho \sum_i \|z_i^{t+1} - z_i^t\|_F^2 \end{aligned}$$

Moreover, from Lem. B.3 we deduce that $\|z_i^{t+1} - z_i^t\|_F^2 \rightarrow 0$ if $t \rightarrow \infty$. Therefore, the global loss converges, $\Phi(z^{t+1}) \rightarrow \hat{\Phi}$ to some constant $\hat{\Phi}$ \square

So far, we only know that our algorithm generates a convergent sequence. It remains to show that the sequence converges to a Boolean solution, which follows in Thm. B.2.

Theorem B.2 (Boolean Convergence of Alg. 1 (restated)). *If λ^t is a monotonically increasing sequence with $\lambda^{t-1} \leq \lambda^t$, $\lim \lambda^t \rightarrow \infty$, and $\lambda^t - \lambda^{t-1} \leq \infty$, then V_1^T, \dots, V_c^T and \widehat{V}^T from the sequence generated by Alg. 1 convergences as $\lim_{T \rightarrow \infty} \text{dist}(\widehat{V}^T, \{0, 1\}) \rightarrow 0$ to a Boolean matrix.*

Proof. In each update round, the client i performs the proximal alternating linear minimization steps laid out in Eq. 7, yielding an updated V_i^t (resp. U_i^t). Focusing on V_i (independent of client-server communication), we first show that the gradient of V_i goes to zero. As shown by Thm. B.1 and Lem. B.3, our sequence of alternating linear optimization steps followed by scaled proximal steps convergence. Note that our gradients are bounded and are Lipschitz continuous. Because we scale our proximal operators with respect the Lipschitz moduli of the respective gradients, notably prevent the proximal operator and gradient steps from alternatingly between 0 and 1, thus creating a convergent sequence to a stable solution. We need to verify that the proximal operator projects to binary solutions, i.e., $\lim_{\lambda^t \rightarrow \infty} \text{prox}(x) \in \{0, 1\}$ for $\lambda^t \rightarrow \infty$. We do this with a case distinction: For $x \leq 0.5$, we obtain $\lim(x - \kappa \text{sign}(x))(1 + \lambda^t)^{-1} = 0$, and analogously for $x > 0.5$, we obtain $\lim(x - \kappa \text{sign}(x - 1) + \lambda^t)(1 + \lambda^t)^{-1} = 1$, thus having ensured a binary proximity, for $\lambda^t \rightarrow \infty$ with $\lambda^t \leq \lambda^{t+1}$ and $\lambda^{t+1} - \lambda^t \leq \infty$, any bounded x , and finite $\kappa \in \mathbb{R}_+$. Therefore, independent of communication rounds, the gradient converges to 0 and the proximal operator converges to a binary solution. It remains to show that for $t \rightarrow \infty$, a binary solution stays stable, meaning that a global binary solution implies local convergence. By assuming that a client in round t receives a binary aggregate \widehat{V} from the server, we obtain $\|\eta \nabla_V \|A_i - U_i^{t-1} V_i\|_{\max} \leq \epsilon$ for $\epsilon < 1/2$. By abbreviating the gradient-step result

$$V' = V_i^{t-1} - \eta \nabla_{V_i} \|A_i - U_i^{t-1} V_i\|_F^2$$

we see that $V'_{pq} < 1/2$ if $[V_i^{t-1}]_{pq} = 0$, and $V'_{pq} > 1/2$ if $[V_i^{t-1}]_{pq} = 1$, which implies that $\text{prox}_{\lambda^t \kappa}(V')$ is binary and $V_i^t = V_i^{t-1}$. Moreover, repeating these steps for \widehat{V}^t , we obtain boolean aggregates upon convergence. \square

Converging Clients

In this part, we demonstrate the convergence of each client in Algorithm 1. Specifically, we show that the decrease between client iterations is sufficiently large, while ensuring convergence to stable solutions. To achieve this, we employ the following lemmas. We establish that the sequence generated by each client converges both in terms of objective function value and to a critical point of the objective function in Lemma B.3. We further provide that the difference of the sequence under finite length conditions is bounded. Subsequently, Lemma B.4 ensures that gradients of the objective function are limited, thereby remain within a certain proximity to the current point. In Lemma B.5, we establish a sufficient decrease property, ensuring that the objective

function decreases at each iteration by a certain amount. By combining these lemmas, we demonstrate the convergence of each client in the algorithm, enabling the global convergence proof in Thm. B.1. In summary, our sub goals are as follows:

1. We aim to demonstrate the convergence of each client.
2. We establish that the decrease between client iterations is sufficiently large.
3. To achieve this, we initially bound all subdifferentials for each client-block, as outlined in Lem. B.4.
4. Subsequently, we utilize this information to bound the gain.

Lemma B.3 (Convergence of client i in Alg. 1). *Let $\{z_i^t \triangleq (U_i^t, V^t)\}_{i \in \mathbb{N}}$ be the sequence generated by a client i in Alg. 1, then*

1. *the client objective $\{\Phi_i(z_i^t)\}_k$ converges to Φ_i^* , and*
2. *the sequence $\{z_i^t\}_k$ converges to a critical point of $\Phi_i(z_i^*)$,*

for $t \rightarrow \infty$, assuming that Φ_i is continuous on $\text{dom } \Phi_i$. Furthermore, if a subsequence z_i^t starts from the shared coefficients \widehat{V} , i.e., $V_i^1 \equiv \widehat{V}$, then the difference $\|V_i^t - \widehat{V}\|_F$ between V_i^t and \widehat{V} is bounded by a finite constant ρ for $t \rightarrow T$.

Before we proof Lem. B.3, we sketch the proof concept as follows. A problem with block-coordinate methods or Gauss-Seidel approaches lies in showing global convergence for these non-convex problems. Attouch, Bolte, and Svaiter (2013) demonstrate the convergence of a sequence generated by a generic algorithm to a critical point of a given proper, lower semicontinuous function Ψ (in our case Φ_i) over a Euclidean space \mathbb{R}^N and establish that the algorithm converges to a critical point of Ψ . To achieve this, we must ensure that the convergence conditions, which are necessary for the convergence of various descent algorithms, are satisfied. If satisfied, they ensure that the set of points of the sequence is nonempty, compact, and connected, with the set being a subset of the critical points of Ψ .

Sufficient Decrease Property This property ensures that with each iteration, the objective value decreases sufficiently. Here the aim is to find a positive constant ρ_1 such that the difference between successive function values decreases sufficiently with each iteration, i.e.,

$$\rho_1 \|z^{t+1} - z^t\|^2 \leq \Psi(z^t) - \Psi(z^{t+1}), \quad \forall t = 0, 1, \dots$$

Subgradient Lower Bound This property ensures that the algorithm does not move too far from the current iterate. Assuming the generated sequence is bounded, we seek another positive constant ρ_2 such that the norm of the difference between consecutive iterates is bounded by a multiple of the norm of the subgradient of Ψ at the current iterate, i.e.,

$$\|w^{t+1}\| \leq \rho_2 \|z^{t+1} - z^t\|, \quad w^t \in \partial \Psi(z^t), \quad \forall t = 0, 1, \dots$$

Because we need a certain stability for our Boolean convergence argument, we have to show that we converge to a critical point. Second, they show global convergence to a critical point using the KL property.

Kurdyka-Łojasiewicz Property To establish global convergence to a critical point, they introduce an additional assumption on the class of functions Ψ being minimized, known as the Kurdyka-Łojasiewicz (KL) property. Intuitively, if this property is satisfied, it prevents the objective to become too flat around a local minimizer, so that the convergence rate would be too low. It does so by creating a locally-convex or linear ‘surrogate’ or ‘gauge’ function g that measures the distance between z and z^*

$$g(\Psi(z) - \Psi(z^*)) \geq \text{dist}(0, \partial\Psi)$$

or more specifically

$$g(\Psi(z) - \Psi(z^*)) \geq \|\partial\Psi\|$$

where, roughly speaking, $z \in \text{Neighborhood}_\eta(z^*)$ (Nesterov 2004). Attouch, Bolte, and Svaiter (2013) have shown that every bounded sequence generated by the proximal regularized Gauss-Seidel scheme converges to a critical point, assuming that the objective function satisfies the KL property (Attouch, Bolte, and Svaiter 2013). We satisfy this assumption, as our objective is comprised of semi-algebraic functions. Now, leveraging the descent property of the algorithm and a uniformization of the KL property, they show that the generated sequence is a Cauchy sequence (Attouch, Bolte, and Svaiter 2013), i.e.,

$$\lim_{l \rightarrow \infty} \sum_{t=l}^{\infty} \|z^t - z^{t-1}\| \rightarrow 0.$$

Proof. (Attouch, Bolte, and Svaiter (2013)). Because Φ_i comprises lower semi-continuous functions on $\text{dom } \Phi_i$, and that all partial gradients are globally Lipschitz, all assumptions for the proof are met (Attouch, Bolte, and Svaiter 2013). Together with (i) sufficient decrease property (Lem. B.5), (ii) lower-bounded subgradients (Lem. B.4), (iii) the Uniformed KL property of (via Lem. B.6), the convergence lemma follows from the global convergence property in Attouch, Bolte, and Svaiter (2013)’s proof. \square

We now formally proof the three properties, i.e., lower-bounded subgradients, sufficient decrease property, and the uniformed KL property.

Lemma B.4 (Lower-bounded Subgradients). *There is a ρ , such that the gradients*

$$\text{dist}(0, \partial\Phi_i(z_i^{t+1})) \leq \rho \|z_i^{t+1} - z_i^t\|_F.$$

Proof. To show that the lemma holds, it suffices that we bound each subgradient in the set $\partial\Phi_i(z_{k+1})$ separately. Focusing on the V_i -block, we want to show

$$\|w^{t+1}\|_F \leq \rho_2 \|V_i^{t+1} - V_i^t\|_F$$

for all $w^{t+1} \in \partial\Phi_i(V_i^{t+1})$ restricted to the V_i^{t+1} -block (analogously repeating the below for U_i^{t+1}). Because the subdifferential of the maximum-term $\max\{r^I(x), r^II(x)\}$ is the union of the subdifferentials of its active parts, and our regularizer is piecewise convex, we obtain three gradients per block:

$$\partial\Phi_i(U_i, V_i) = \nabla_{V_i} \frac{1}{2} \|A_i - U_i V_i\|_F^2 + \nabla_{V_i} \frac{1}{2} \|V_i - V_i^t\|_F^2 + \partial R(V_i),$$

$$\partial R(V_i) = \begin{cases} \nabla_{V_i} r^I(V_i) & R(V_i) < R(V_i - 1) \\ \text{conv}(\nabla_{V_i} r^II(V_i), \nabla_{V_i} r^I(V_i)) & R(V_i) = R(V_i - 1) \\ \nabla_{V_i} r^II(V_i) & R(V_i) > R(V_i - 1) \end{cases}.$$

Next, we bound the norm of the first subdifferential

$$\begin{aligned} & \|\nabla_{V_i} \Phi_i + \partial R(V) + \frac{\gamma}{2} (V - V_i^t)\|_F \\ & \leq \|\nabla_{V_i} \Phi_i + \partial R(V_i)\|_F + \frac{\gamma}{2} \|V - V_i^t\|_F \\ & \leq \frac{\rho}{2} \|V - V_i^t\|_F + \frac{\gamma}{2} \|V - V_i^t\|_F \\ & \leq \max\{\rho, \gamma\} \frac{1}{2} \|V - V_i^t\|_F. \end{aligned}$$

Repeating for the other cases, the total bound ρ is the maximum per block and per subdifferential bounds. Based on Lem. B.5, under the assumption that Φ_i is continuous on its domain, and provided that there exists a convergent subsequence (i.e., condition (a)), the continuity condition required in (Attouch, Bolte, and Svaiter 2013) holds, i.e., there exists a subsequence $\{z_i^t\}_{k \in \mathbb{N}}$ and a point z_i^* such that

$$z_i^t \rightarrow z_i^* \quad \text{and} \quad \Phi_i(z_i^t) \rightarrow \Phi_i(z_i^*) \quad \text{as } t \rightarrow \infty.$$

\square

Lemma B.5 (Sufficient Decrease Property). *For the sequence of points $\{z^t\}_k$ generated by the block-coordinate method in Alg. 1, then*

$$\Phi_i(z_i^{t+1}) \leq \Phi_i(z_i^t) - \rho_1 \|z_i^{t+1} - z_i^t\|_F^2.$$

Proof. The loss function for the V_i -block in our local block-coordinate descent is

$$\|A_i - U_i^t V_i^{t+1}\|_F^2 + \|V_i^{t+1} - \widehat{V}\|_F^2 + R(V_i^{t+1}).$$

Likewise for the U_i -block

$$\|A_i - U_i^{t+1} V_i^{t+1}\|_F^2 + R(U_i^{t+1}).$$

After taking a gradient step, Alg. 1 proceeds with a *Boolean projection* regarding R (for U and V blocks) and a *proximity projection* to \widehat{V} (only for V).

We proceed with the V_i -block, while the proof for the U_i -block is analogous. First, the **Boolean proximal** projection operator $\text{prox}_R(V^t)$ yields a *minimizer* of the optimization problem

$$\overset{I}{V}^k \leftarrow \arg \min_Y \frac{1}{2} \|V^t - Y\|_F^2 + R(Y).$$

By definition, $\overset{I}{V}_i^k$ lies in a ρ_I -bounded proximity to V_i^t . Second, the **proximity proximal** projection operator $\text{prox}_{\gamma \widehat{V}}^P(\overset{I}{V}_i^k)$ is the minimizer of

$$\overset{II}{V}_i^k \leftarrow \arg \min_Y \frac{1}{2} \|\overset{I}{V}_i^k - Y\|_F^2 + \nu^{1/2} \|\widehat{V} - Y\|_F^2.$$

By definition, $\overset{II}{V}_i^k$ lies in the ρ_{II} -bounded proximity to $\overset{I}{V}_i^k$. Repeating for the U_i -blocks and using a transitivity argument,

by using that our gradients have finite Lipschitz moduli, we conclude that both projections lie in a ρ -bounded region around z_i^t .

Using the following relationships,

$$\begin{aligned}\Phi_i(z_i^{t+1}) &\leq \Phi_i(z_i^t) + \rho \|z_i^{t+1} - z_i^t\|_{\mathbb{F}}^2, \\ \Phi_i(z_i^{t+1}) - \rho \|z_i^{t+1} - z_i^t\|_{\mathbb{F}}^2 &\leq \Phi_i(z_i^{t+1}), \text{ and} \\ \Phi_i(z_i^{t+1}) &\leq \Phi_i(z_i^t).\end{aligned}$$

we now bound the loss reduction in terms of the norm of differences in the following.

$$\begin{aligned}\Phi_i(z_i^{t+1}) &\leq \Phi_i(z_i^t) \\ \Phi_i(z_i^{t+1}) - \rho \|z_i^{t+1} - z_i^t\|_{\mathbb{F}}^2 &\leq \Phi_i(z_i^t) \\ \Phi_i(z_i^{t+1}) - \rho \|z_i^{t+1} - z_i^t\|_{\mathbb{F}}^2 &\leq \Phi_i(z_i^{t-1}) + \rho \|z_i^t - z_i^{t-1}\|_{\mathbb{F}}^2 \\ \Phi_i(z_i^{t+1}) - \rho \|z_i^{t+1} - z_i^t\|_{\mathbb{F}}^2 &\leq \Phi_i(z_i^t) + \rho \|z_i^t - z_i^{t-1}\|_{\mathbb{F}}^2 \\ \Phi_i(z_i^{t+1}) - \Phi_i(z_i^t) &\leq \rho \|z_i^t - z_i^{t-1}\|_{\mathbb{F}}^2 + \rho \|z_i^{t+1} - z_i^t\|_{\mathbb{F}}^2 \\ \Phi_i(z_i^{t+1}) - \Phi_i(z_i^t) &\leq \rho \|z_i^{t+1} - z_i^t\|_{\mathbb{F}}^2\end{aligned}$$

□

Lemma B.6 (Uniformized Kurdyka-Łojasiewicz (KL)). Φ_i is a KL function.

Proof. Φ_i function is composed of p -norms ($p \in \{1, 2\}$), and indicator functions, and therefore satisfy the KL-property (At-touch, Bolte, and Svaiter 2013). □

C Competitors

For a given aggregation function (such as rounded averaging (10), majority voting (11), or logical $\circ r$ (12)), we summarize the federation strategy of centralized BMF algorithms in Alg. 2.

Algorithm 2: Aggregated BMF

Input: C clients with local matrices A_1, \dots, A_C , local BMF algorithm \mathcal{A} , aggregation function aggregate

Output: local feature matrices U_1, \dots, U_C , global coefficient matrix \widehat{V}

- 1 **Locally at client i do**
 - 2 $U_i, V_i \leftarrow \mathcal{A}(A_i)$
 - 3 **Centrally at server do**
 - 4 receive V_1, \dots, V_C
 - 5 $\widehat{V} \leftarrow \text{aggregate}(V_1, \dots, V_C)$
 - 6 transmit \widehat{V} to all clients
 - 7 **Locally at client i do**
 - 8 receive \widehat{V} from the server
 - 9 assign $V_i \leftarrow \widehat{V}$
-

Obtaining Boolean Matrices from ZHANG’s Factorization

The relaxation-based binary matrix factorization of ZHANG (Zhang et al. 2007) does not necessarily yield Boolean factors

Table 2: Real-world datasets from 4 diverse domains. We show extents, density, and the selected number of components for 10 real-world datasets.

Dataset	Rows	Cols	Density	Components
ACS Inc	1630167	998	0.010	20
ACS Pov	3271346	836	0.024	20
cs.LG	145981	14570	0.005	50
Goodreads	350332	9374	0.001	50
HPA Brains	76533	20082	0.239	100
Movielens	162541	62423	0.002	20
Netflix	480189	17770	0.007	20
TCGA	10459	20530	0.019	33

upon convergence. Furthermore, this method yields matrices that do not lend themselves to rounding, such that in practice, rounding does not yield desirable results *unless* the rounding threshold is carefully chosen. To choose well-factorizing rounding thresholds, we take inspiration from PRIMP (Hess, Morik, and Piatkowski 2017), searching those thresholds that minimize the reconstruction loss,

$$\sum_{c \in [C]} \|A^c - [U_{ij}^c \geq \alpha]_{ij} \circ [V_{ij}^c \geq \beta]_{ij}\|,$$

from the equi-distant grid between 1×10^{-12} and 1 containing 100 points in each direction.

D Datasets

To explore the realm of **recommendation systems**, we have included *Goodreads* (Kotkov et al. 2022) for book recommendations, as well as *Movielens* (Harper and Konstan 2015) and *Netflix* (Netflix, Inc. 2009) for movie recommendations. To focus on positive ratings, we binarized user ratings, setting ratings ≥ 3.5 to 1.

In the field of **life sciences**, we consider cancer genomics through *TCGA* (Institute 2005) and single-cell proteomics using *HPA* (Bakken et al. 2021; Sjöstedt, Zhong, and et. al 2020). Specifically, *TCGA* records 1s for gene expressions in the upper 95% quantile and *HPA* records by 1 if RNA has been observed in single cells.

For **social science** inquiries, we investigate poverty (P) and income (I) analysis using the *Census* (U.S. Census Bureau 2023) dataset. To binarize, we employ one-hot encoding based on the features recommended by Folktables (Ding et al. 2021).

In the domain of **natural language processing**, we focus on higher-order word co-occurrences using ArXiv abstracts from the cs.LG category (Collaboration 2023). Each paper corresponds to a row whose columns are 1 if the corresponding word in our vocabulary has been used in its abstract. The vocabulary consists of words with a minimum frequency of 1 ‰ in ArXiv cs.LG abstracts (*cs.LG R*) and their lemmatized, stop-word-free counterparts (*cs.LG*).

We summarize extents, density, and chosen component counts for each real-world dataset in Appendix E, Table 2.

E Reproducibility

Supplementing the information provided in Sec. 4, here, we provide hyperparameter choices for FELB and FELB^{MU}. We use the iPALM optimization approach for FELB and FELB^{MU}. Because both algorithms exhibited relatively stable performance fluctuations when it came to tuning, we used the same set of hyperparameters for each experiment and each dataset, thus omitting the commonly necessary hyperparameter tuning step. In all experiments with FELB and FELB^{MU}, we used the regularizer coefficients $\lambda = 0.1$ and $\kappa = 0.001$, a regularization rate $\lambda_t = \lambda \cdot 1.05^t$, an iPALM inertial parameter $\beta = 0.001$, a maximum number of iterations of 100, and a number of local rounds per iteration of 1, 10, or 50, as indicated by the experiments. For ELBMF, we choose $\kappa = 0.01$, $\lambda = 0.01$, $\lambda_t = \lambda \cdot 1.02^t$, and $\beta = 0.01$. We provide ZHANG and ELBMF with a larger iteration limit of 1 000, multiplying FELB’s local rounds by its iteration count. For ASSO, we set gain, loss, and threshold parameters to 1.0. For MEBF, we use a threshold of 0.5 and a cover limit of 0.95.

F Additional Experiments

Complementing the discussion in Sec. 4, here, we show additional results for ASSO, GRECOND, MEBF, ELBMF, and ZHANG, as well as FELB and FELB^{MU}, for all experiments. We focus on the quantification not present in the main body of this paper. Here, we aim to answer the following additional questions.

- Q4 How does client drift impact real-world performance?
- Q5 How different the post-hoc aggregations for BMF are?
- Q6 How stably does our methods converge?
- Q7 How robust do we handle client drift?
- Q8 How achievable is differential privacy in different circumstances?

Real-world Experiments

In addition to results presented in Table 1, we provide the RMSD in Table 3, where we see that the FELB^{MU} and FELB are the two best performing methods, followed by GRECOND.

Real-world Drift Experiments

Next, because the performance depends on the communication frequency, we evaluate our method in 3 different scenarios: Rare (max 50 client epochs), Occasional (max 10 epochs), and Frequent synchronizations (every round). To visualize relative performance differences, we compute the *relative RMSD*

$$\frac{\text{RMSD}(\text{FELB})}{\text{RMSD}(\text{FELB}^{\text{MU}})},$$

depicted in Fig. 5 for all real-world datasets in different synchronization regimes. Because ASSO, GRECOND, MEBF, ZHANG, and ELBMF are not directly federated, they are independent of the change in communication frequency and therefore omitted. In Fig. 5, we see that our algorithm maintain a high prediction performance regardless of the communication overhead. We observe that FELB and FELB^{MU} perform similarly well under occasional and frequent communications. We observe a shrinking performance gap between FELB and

FELB^{MU} when increasing the communication frequency, almost reaching the same performance. This indicates that FELB’s larger gradient-step-sizes are responsible for a higher client drift, which is mitigated by a high communication frequency. Regardless of being under occasional and frequent communication regime, FELB and FELB^{MU} are the highest performing algorithms.

Post-hoc Aggregations

As there is no prior art specifically for aggregation federated BMF clients, we seek experimentally answer which of the equations Eqs. (10)–(11) yield the lowest reconstruction loss. To this end, we consider a growing number of synthetic abundant data as described for Fig. 3. While we observe in Fig. 6 and in Fig. 7 that *rounded average* and *consensus voting* are performing similarly, both significantly outperform *logical or*. For brevity, we therefore mostly report results for *consensus voting* in Sec. 4.

Empirical Convergence

This study aims to investigate the empirical convergence properties of the proposed methods. In this study, we examine the empirical convergence properties of our methods. We generate synthetic data according to the procedure outlined in Sec. 4. We then measure the reconstruction loss as the number of global iteration steps increases. Fig. 8 demonstrates that our methods rapidly converge to a lower loss corresponding to non-Boolean solutions. Following a swift initial decrease, the loss only minimally increases as we approach a feasible Boolean solution upon convergence.

Client Drift

We aim to understand the impact of infrequent synchronizations on the convergence results. To investigate this, we vary the number of local iterations per client from 1 (frequent synchronizations) to 50 (infrequent synchronizations), using synthetic data. In Fig. 9, we observe that the loss is significantly affected by the increasing number of iterations. We see that the loss flattens-out after approximately 25 client local optimization epochs before synchronization. While our methods achieve a reasonably high F_1^* -score with respect to the ground-truth—even with infrequent synchronizations—our competitors do not show similar results.

Differential Privacy

We aim to understand how differential privacy impacts reconstruction quality. Previously, we studied the effect of clipped noise mechanisms (Fig. 4). Here, we extend this experiment to include non-clipped noise mechanisms, as shown in Figures 10 and 11. Specifically, we apply non-clipped Gaussian and Laplacian noise to federated factorization algorithms that operate on real-valued numbers, while limiting discrete Boolean factorization algorithms to Bernoulli noise.

In Fig. 10, we observe that the F_1 -score decrease significantly only at high differential privacy coefficients. At moderate levels, we achieve differentially private reconstructions using both clipped and non-clipped Gaussian and Laplacian noise mechanisms, as well as Bernoulli ‘XOR’ noise. In

Table 3: FELB and FELB^{MU} consistently achieve top performances. We illustrate the RMSD of ASSO, GRECOND, MEBF, ELBMF, and ZHANG under voting aggregation, as well as federated FELB, and FELB^{MU} on 8 real-world data across 50 clients. We highlight the best algorithm with **bold**, the second best with underline, and indicate missing data by a dash –.

Dataset	ASSO ^V	MEBF ^V	GRECOND ^V	ZHANG ^V	ELBMF ^V	FELB ^{MU}	FELB
ACS Inc	4.583	4.929	3.485	5.005	5.005	<u>3.962</u>	4.560
ACS Pov	<u>5.822</u>	–	4.785	7.734	7.190	7.576	7.588
cs.LG	–	3.398	3.350	3.398	3.398	<u>3.372</u>	3.396
Goodreads	–	1.669	1.668	–	1.669	1.641	<u>1.660</u>
HPA Brain	–	19.537	–	24.434	24.434	<u>24.409</u>	24.433
MovieLens	–	1.956	–	–	1.962	1.914	<u>1.925</u>
Netflix	–	4.075	–	–	4.084	3.982	<u>4.009</u>
TCGA	6.858	6.834	6.871	6.872	6.872	6.346	<u>6.420</u>
<i>Rank</i>	4.000	3.625	2.875	5.000	4.625	1.750	<u>2.750</u>

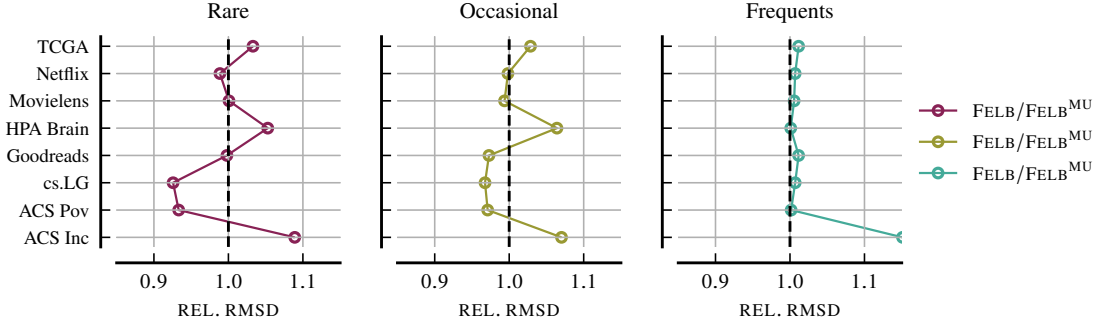


Figure 5: FELB and FELB^{MU} perform similarly when we synchronize clients frequently, while FELB^{MU} tends to improve over FELB under low-communication regimes. We show the RMSD-fraction (F) of FELB and FELB^{MU} as *relative* score on real-world datasets with varying communication frequencies.

Fig. 11, we see that the reconstruction loss follows a similar trend for both Gaussian and Laplacian noise mechanisms. The Bernoulli mechanism, however, results in a much lower reduction in RMSD than in the F_1 . Although all methods exhibit a similar trend, FELB and FELB^{MU} demonstrate robustness regarding differential privacy, consistently outperforming competitors in terms of RMSD and F_1 score.

G Limitations

Our research is motivated by learning from private Boolean data generated by similar sources, situated at few research centers. As such, we focus on suitable experiments in our research, while we abstain from distant but related problems.

Firstly, our approach does not incorporate personalized federated learning (PFL), which could potentially enhance individual client performance by tailoring the model to specific client data. Additionally, our experimental study does not address heterogeneous data distributions across clients, which is a common scenario in real-world applications. Furthermore, our focus is on learning and knowledge discovery from federations involving a limited number of clients, specifically in the context of research centers. This is in contrast to scenarios involving millions of clients, such as those sometimes encountered in different federated learning applications.

We experimentally demonstrate under which circumstances our method breaks, involving experiments with noise levels 4, privacy levels 4, client counts 4, dataset sizes 4, client-server communication intervals F , and dataset domains D , thereby providing an extensive overview over practical strength and weaknesses.

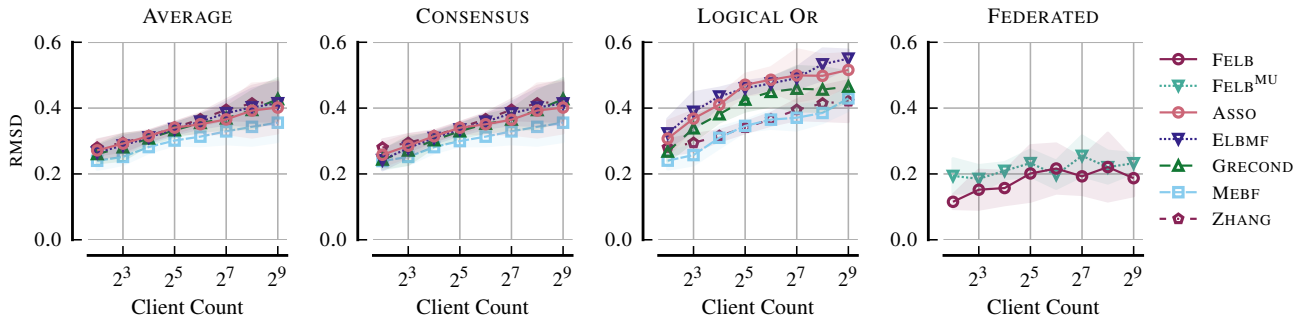


Figure 6: The Boolean matrix aggregation methods *rounded average* and *consensus voting* significantly outperform *logical or*. We show the loss for post-hoc aggregated BMF methods, for growing client count with synthetic abundant data.

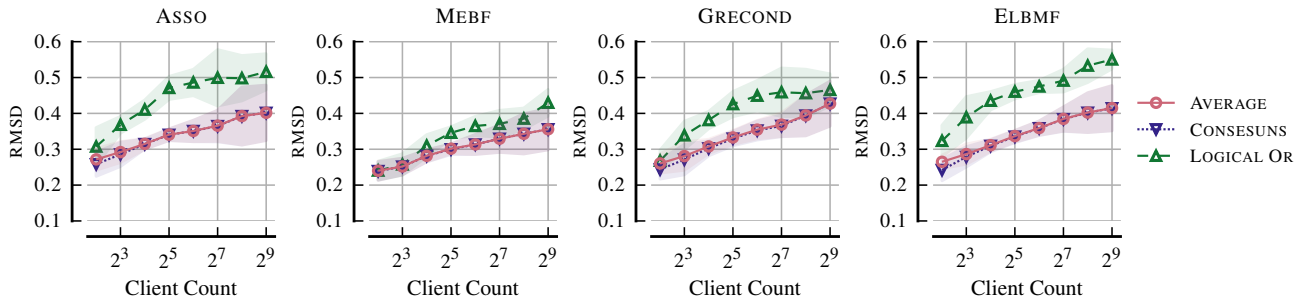


Figure 7: The Boolean matrix aggregation methods *rounded average* and *consensus voting* significantly outperform *logical or*, depicting results specifically for post-hoc aggregated BMF methods, for growing client count with synthetic abundant data.

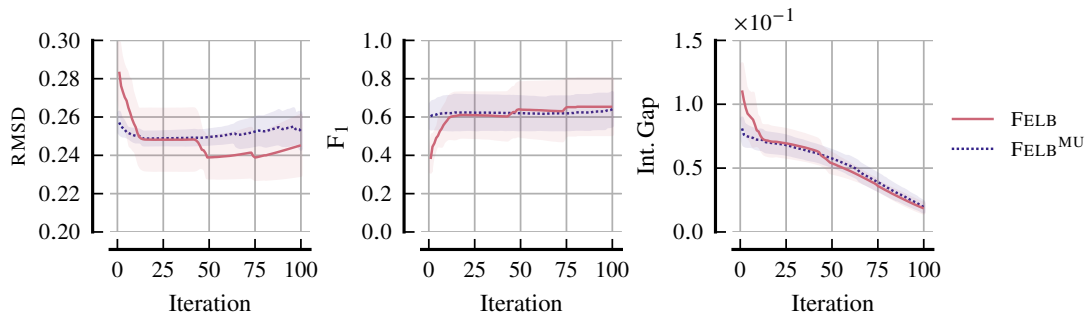


Figure 8: Our methods rapidly achieve a lower reconstruction loss for non-Boolean solutions and maintain minimal loss increase while approaching a feasible Boolean solution. We illustrate the history of loss, F₁ score, and integrality gap over increasing number of iterations.

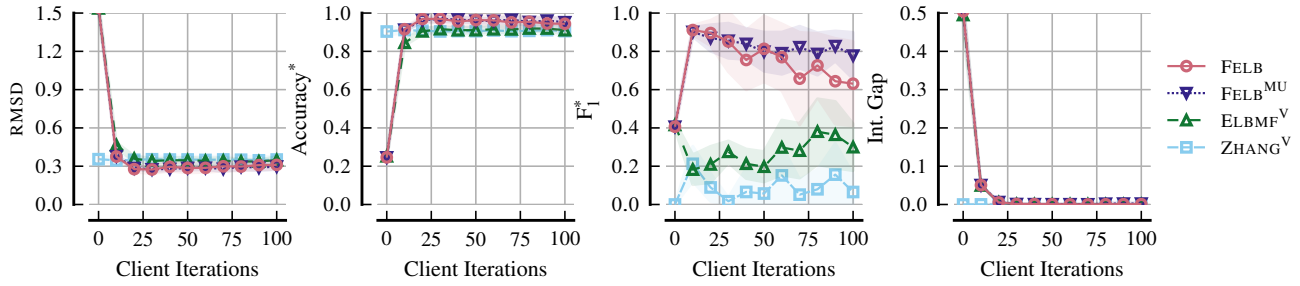


Figure 9: Our algorithm demonstrates robustness in achieving high convergence rates despite infrequent synchronizations. We illustrate the history of loss, F₁ score, F₁^{*} score regarding ground-truth, and integrality gap over increasing number of local per-client iterations before global synchronizations.

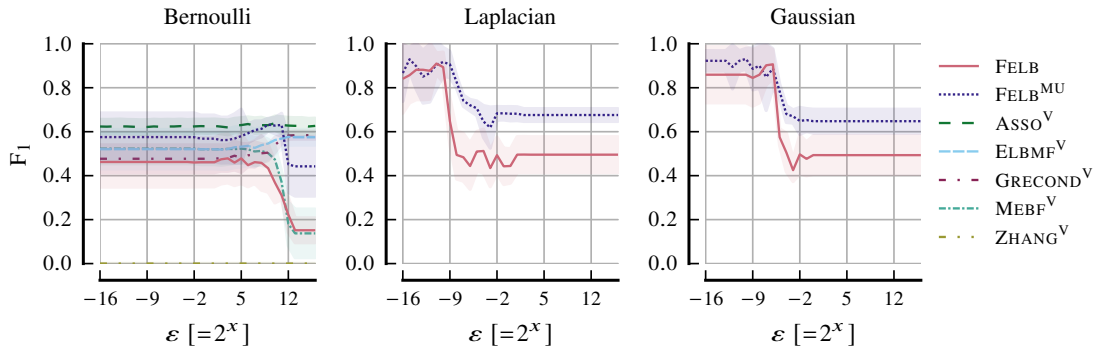


Figure 10: Our algorithms largely maintains the prediction performance for moderately high differential privacy coefficients. We depict the F₁-score trend across various levels of differential privacy, for non-clipped Gaussian and Laplacian noise mechanisms, as well as the Bernoulli ‘XOR’ noise mechanism.

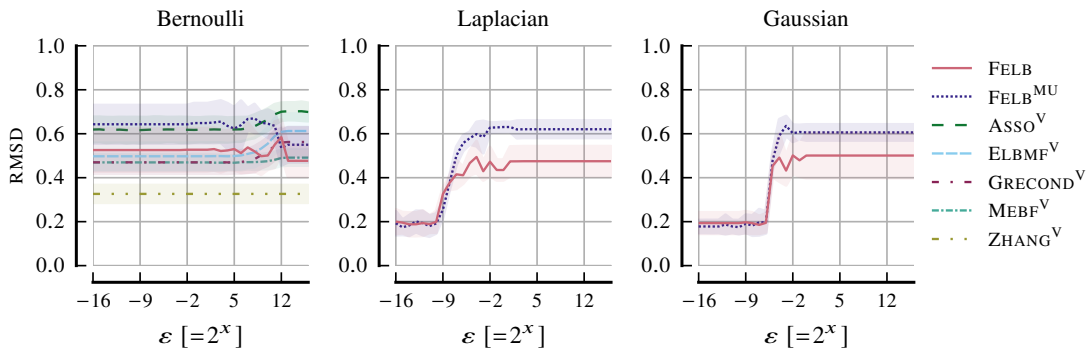


Figure 11: Our algorithms largely maintains the reconstruction quality for moderately high differential privacy coefficients. We depict the reconstruction loss trend across various levels of differential privacy, for non-clipped Gaussian and Laplacian noise mechanisms, as well as the Bernoulli ‘XOR’ noise mechanism.