

Supervised Matrix Factorization: Local Landscape Analysis and Applications

Joowon Lee,¹ Hanbaek Lyu,² and Weixin Yao³

1: Department of Statistics, University of Wisconsin – Madison, WI, USA
 2: Department of Mathematics, University of Wisconsin – Madison, WI, USA
 3: Department of Statistics, University of California, Riverside, CA USA



Abstract

Supervised matrix factorization (SMF) is a classical machine learning method that seeks low-dimensional feature extraction and classification tasks at the same time. Training an SMF model involves solving a non-convex and factor-wise constrained optimization problem with at least three blocks of parameters. Due to the high non-convexity and constraints, theoretical understanding of the optimization landscape of SMF has been limited. In this paper, we provide an extensive local landscape analysis for SMF and derive several theoretical and practical applications. Analyzing diagonal blocks of the Hessian naturally leads to a block coordinate descent (BCD) algorithm with adaptive step sizes. We provide global convergence and iteration complexity guarantees for this algorithm. Full Hessian analysis gives minimum L^2 -regularization to guarantee local strong convexity and robustness of parameters. We establish a local estimation guarantee under a statistical SMF model. We also propose a novel GPU-friendly neural implementation of the BCD algorithm and validate our theoretical findings through numerical experiments. Our work contributes to a deeper understanding of SMF optimization, offering insights into the optimization landscape and providing practical solutions to enhance its performance.

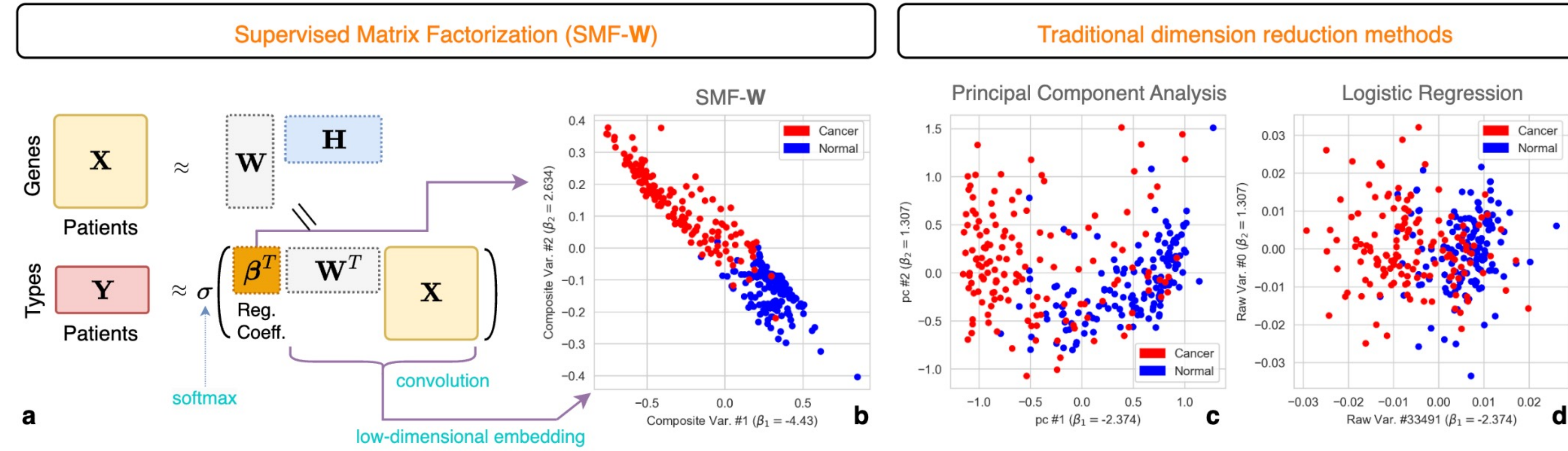


Figure 1. (a) Overall scheme of Supervised Matrix Factorization (specifically, SMF- \mathbf{W} with rank $r = 2$). The columns of \mathbf{W} serve as ‘composite variables’ or ‘filters’, whose association with the labels is given by the regression coefficients in β . Taking convolution of the raw data matrix \mathbf{W} with \mathbf{W} gives a supervised dimension reduction, as illustrated in (b) for a 35,982-dimensional gene microarray data for breast cancer patients. Similar dimension reduction results obtained by (c) principal component analysis along with logistic regression and (d) logistic regression to select the two most highly associated raw variables show less clear separation.

1. Objective and Model Formulation

- Labeled data (y_i, \mathbf{x}_i) , $i = 1, \dots, n$:
- y_i = binary label (e.g., 1=“Cancer”, 0=“Normal”)
- \mathbf{x}_i = high-dimensional feature vector (e.g., gene expression data)
- Dimension reduction + Classification** at the same time?
- Key difficulty:

2. How the model works

- The filter matrix \mathbf{W} is learned to serve double purposes:
 - Input for Classification : $\mathbf{W}^T \mathbf{x}_i = r$ – dim compressed feature
 - Feature reconstruction : $\mathbf{X} \approx \mathbf{W}\mathbf{H}$ for some $\mathbf{H} \in \mathbb{R}_{\geq 0}^{r \times n}$
- Once the filter \mathbf{W} and reg. coefficients β are learned:

$$\mathbb{P}(y_i = 1 \mid \mathbf{x}_i) = \frac{\exp(\mathbf{a}_i)}{1 + \exp(\mathbf{a}_i)}, \quad \text{where } \mathbf{a}_i = \beta^T \mathbf{W}^T \mathbf{x}_i$$

3. Local Optimization Landscape: Diagonal

- To sketch the idea, consider the matrix factorization loss only:

$$f(\mathbf{W}, \mathbf{H}) = \|\mathbf{X} - \mathbf{W}\mathbf{H}\|_F^2$$

$$\nabla^2 f = \begin{bmatrix} \text{vec}(\mathbf{W}) & \text{vec}(\mathbf{H})^T \\ \text{vec}(\mathbf{H}) & \begin{bmatrix} \mathbf{H}\mathbf{H}^T \otimes \mathbf{I}_p & \mathbf{A}_{12} \\ \mathbf{A}_{12}^T & \mathbf{I}_n \otimes \mathbf{W}^T \mathbf{W} \end{bmatrix} \end{bmatrix}$$

$\lambda_{\max}(\nabla^2 f)$ Unbounded $\odot \rightarrow$ PGD very sensitive on step size

$$\lambda_{\max}(\nabla_{\mathbf{W}}^2 f(\cdot, \mathbf{H})) = \lambda_{\max}(\mathbf{H}\mathbf{H}^T) \quad \text{Bounded for fixed } \mathbf{H} !! \odot$$

$$\lambda_{\max}(\nabla_{\mathbf{H}}^2 f(\mathbf{W}, \cdot)) = \lambda_{\max}(\mathbf{W}^T \mathbf{W}) \quad \text{Bounded for fixed } \mathbf{W} !! \odot$$

Adaptive BCD:

$$\mathbf{W} \leftarrow \Pi \left(\mathbf{W} - \frac{1}{\lambda_{\max}(\mathbf{H}\mathbf{H}^T) + \varepsilon} (\mathbf{W}\mathbf{H} - \mathbf{X})\mathbf{H}^T \right),$$

$$\mathbf{H} \leftarrow \Pi' \left(\mathbf{H} - \frac{1}{\lambda_{\max}(\mathbf{W}\mathbf{W}^T) + \varepsilon} \mathbf{W}^T (\mathbf{W}\mathbf{H} - \mathbf{X}) \right)$$

Largest Eval of diagonal blocks of the Hessian

<< Largest Eval of the entire Hessian

- We derive Adaptive BCD for SMF similarly
- $$\begin{bmatrix} \text{vec}(\mathbf{W})^T & \text{vec}(\mathbf{H})^T & \text{vec}(\beta)^T & \text{vec}(\Gamma)^T \\ \text{vec}(\mathbf{W}) & \begin{bmatrix} A_{11} & A_{12} & A_{13} & \mathbf{0} \\ A_{21} & A_{22} & \mathbf{0} & \mathbf{0} \\ A_{31} & \mathbf{0} & A_{33} & A_{34} \\ \mathbf{0} & \mathbf{0} & A_{43} & A_{44} \end{bmatrix} \end{bmatrix}$$

4. Local Optimization Landscape: Off-diagonal

$$\|\mathbf{X} - \mathbf{W}\mathbf{H}\|_F^2 + \lambda_1 \|\mathbf{W}\|_F^2 + \lambda_2 \|\mathbf{H}\|_F^2$$

- Smallest L_2 -reg. to ensure *local strong convexity*? (\rightarrow Robust parameter estimation)
- Smallest L_2 -reg. for block diagonal dominance in $\nabla^2 f$

$$\lambda_{\min}(\mathbf{H}\mathbf{H}^T) + \lambda_1 - \|\mathbf{A}_{12}\|_2 > 0$$

$$\lambda_{\min}(\mathbf{W}^T \mathbf{W}) + \lambda_2 - \|\mathbf{A}_{12}\|_2 > 0$$
- High-dim regime: $p \gg n \rightarrow$ No need to reg. $\mathbf{H}!!$

$$\lambda_{\min}(\mathbf{W}^T \mathbf{W}) = \Theta(rp) \gg \Theta(r\sqrt{pn}) = \|\mathbf{A}_{12}\|_2$$

$$\rightarrow \mathbf{H}^{opt} \text{ can be recovered robustly despite non-convexity}$$
- Large-sample regime: $p \ll n \rightarrow$ No need to reg. $\mathbf{W}!!$

$$\lambda_{\min}(\mathbf{H}\mathbf{H}^T) = \Theta(rn) \gg \Theta(r\sqrt{pn}) = \|\mathbf{A}_{12}\|_2$$

Theorem 4.5. (Regularized local consistency) Consider the generative SMF- \mathbf{W} model (16). Assume that Assumptions 4.1 and 4.2 hold. Suppose $\rho := \min_{1 \leq i \leq 4} (\lambda_i - \lambda_{i*}) > 0$.

Suppose $\Lambda_1 > 0$, $\lambda_1 = 0$, and $\sigma \ll 1$ (resp., $\Lambda_2 > 0$ and $\lambda_2 = 0$). Fix $\varepsilon > 0$. Then there exists a constant $C > 0$ such that with probability at least $1 - \varepsilon$, \mathcal{L} in (17) is minimized locally at some $(\hat{\mathbf{W}}, \hat{\theta}, \hat{\lambda})$ (resp., $(\hat{\mathbf{H}}, \hat{\theta}, \hat{\lambda})$) with

$$\|\hat{\mathbf{W}} - \mathbf{W}_*\| \leq C/\sqrt{n} \quad (\text{resp., } \|\hat{\mathbf{H}} - \mathbf{H}_*\| \leq C/\sqrt{n}) \quad (18)$$

$$\|\hat{\lambda} - \lambda_*\| \leq C/\sqrt{n}$$

$$\|\hat{\theta} - \theta_*\|_F \leq Cn^{-1/2} \left(1 + \frac{3 \max\{\lambda_2, \lambda_3, \lambda_4\}}{\rho} \|\theta_*\|_F \right),$$

where $\theta' := (\mathbf{H}', \beta', \Gamma')$, $\theta_* := (\mathbf{H}_*, \beta_*, \Gamma_*)$ (resp., $\theta' := (\mathbf{W}', \beta', \Gamma')$, $\theta_* := (\mathbf{W}_*, \beta_*, \Gamma_*)$) and $\|\theta_*\|_F$ is assumed to be sufficiently small.

5. Experiments

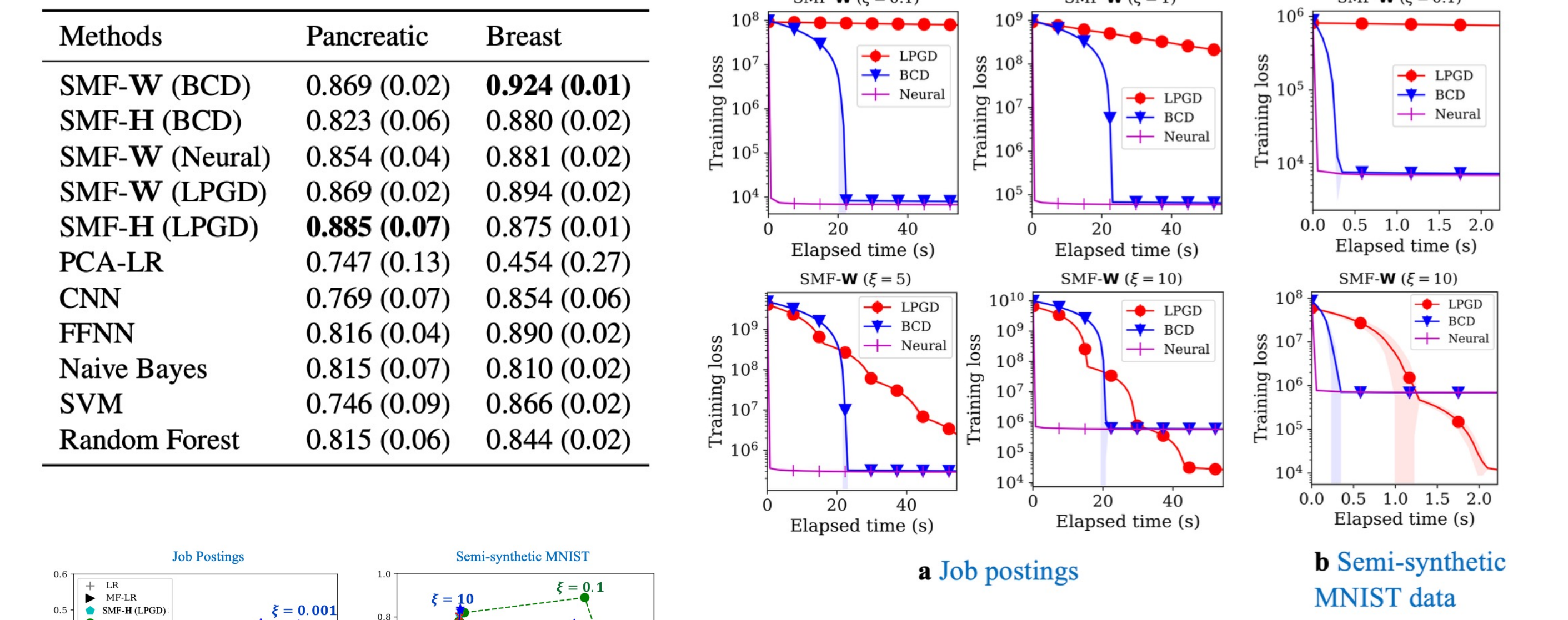
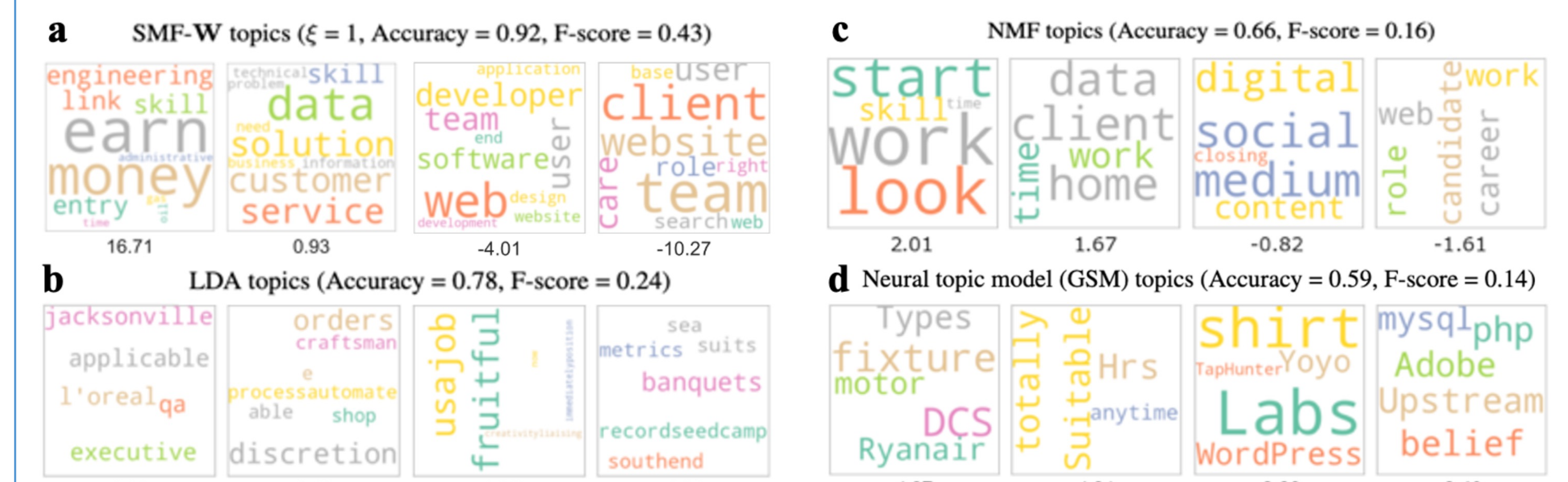


Figure 3. Plots of training loss vs. elapsed time at different ξ values for fitting SMF- \mathbf{W} using Algorithm 1 (BCD), the neural implementation in Figure 2 (Neural), and low-rank projected gradient descent (LPGD) in (Lee et al., 2023). Shaded regions indicate one standard deviation across 10 runs.



6. References

Lee, J., Lyu, H., and Yao, W. Exponentially convergent algorithms for supervised matrix factorization. *NeurIPS 2023*

7. Acknowledgements

NSF DMS-2206296, DMS—010035, DMS-2210272