

# End-to-end Learning, with or without Labels

Corinne Jones, Vincent Roulet, Zaid Harchaoui

University of Washington  
 {cjones6, vroulet, zaid}@uw.edu

## Abstract

We develop a framework **unifying unsupervised and supervised** end-to-end learning. The objective function adjusts gracefully to the amount of supervision, reducing to a **clustering** objective when only unlabeled data is available and to a **classification** objective when only labeled data is available.

## X-Supervised Objective

Consider observations  $x_i, i = 1, \dots, n$ , whose labels  $y_i^* \in \{0, 1\}^k$  may or may not be observed. Given a deep network  $\phi$ , we optimize over the unknown labels and the parameters  $V$  of the network and  $W, b$  of a classifier:

$$\min_{Y \in \mathcal{C}, V, W, b} \frac{1}{n} \sum_{i=1}^n \|y_i - W^T \phi(x_i; V) - b\|_2^2 + \alpha \sum_{j=1}^m \|V_j\|_F^2 + \lambda \|W\|_F^2 - \rho \sum_{i=1}^n \|\phi(x_i; V) - \bar{\phi}\|_2^2$$

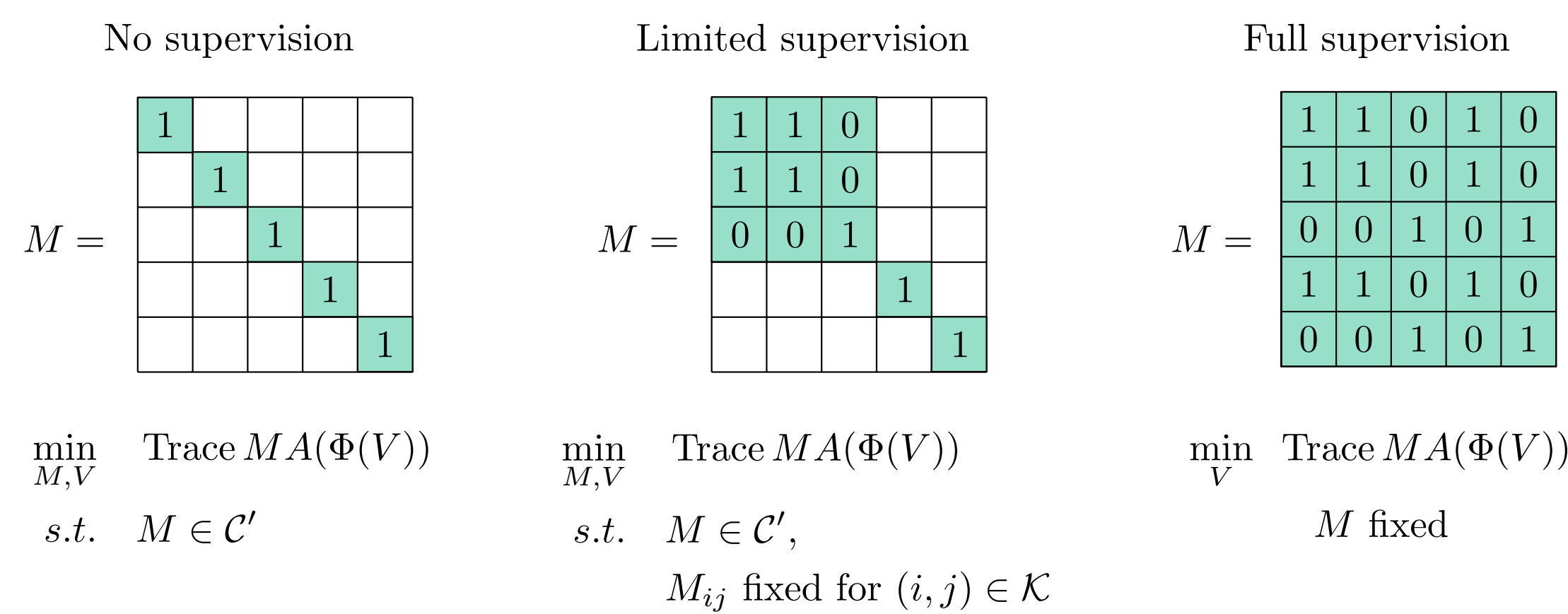
where

$$\mathcal{C}' := \{Y \in \{0, 1\}^{n \times k} : Y \mathbf{1}_k = \mathbf{1}_n, y_i = y_i^* \text{ for } i \in \mathcal{S}, n_{\min} \mathbf{1}_k \leq Y^T \mathbf{1}_n \leq n_{\max} \mathbf{1}_k\}$$

and  $\bar{\phi} = 1/n \sum_{i=1}^n \phi(x_i; V)$ . This extends the objective of Bach and Harchaoui (2007) to the deep setting.

The additional penalty and constraints are interpreted as follows:

- $\rho \sum_{i=1}^n \|\phi(x_i; V) - \bar{\phi}\|_2^2$ : Prevents mapping all observations to the same embedding, i.e.,  $\phi(x_1; V) = \phi(x_2; V) = \dots = \phi(x_n; V)$
- $n_{\min} \mathbf{1}_k \leq Y^T \mathbf{1}_n \leq n_{\max} \mathbf{1}_k$ : Prevents assigning all observations to the same cluster, i.e.,  $y_1 = y_2 = \dots = y_n$



Example equivalence matrix  $M = YY^T$  and problem for varying levels of supervision. For simplicity,  $\alpha = \rho = 0$ .

After optimizing over  $W$  and  $b$  in closed form, we obtain the problem

$$\min_{M, V} \lambda \text{tr}[MA(\Phi(V))] + \alpha \sum_{j=1}^m \|V_j\|_F^2 - \rho \sum_{i=1}^n \|\phi_i(V) - \bar{\phi}\|_2^2 \quad (1)$$

s.t.  $M_{ij} = m_{ij} \quad \forall (i, j) \in \mathcal{K}$   
 $n_{\min} \mathbf{1}_n \leq M \mathbf{1}_n \leq n_{\max} \mathbf{1}_n$   
 $n_{\min} \mathbf{1}_n \leq M^T \mathbf{1}_n \leq n_{\max} \mathbf{1}_n$

where

- $M = YY^T$
- $\phi_i(V) = \phi(x_i; V)$  and  $\Phi(V) = (\phi_1(V), \dots, \phi_n(V))^T$ .
- $A(\Phi) = \Pi_n (\Pi_n \Phi \Phi^T \Pi_n + n\lambda \mathbf{I})^{-1} \Pi_n$  and  $\Pi_n = \mathbf{I}_n - \mathbf{1}_n \mathbf{1}_n^T / n$
- $\mathcal{K}$  is the set of indices of known entries  $m_{ij}$  of  $M := YY^T$ .

## Optimization

We optimize over  $V$  by repeating the following steps on mini-batches:

1. Obtain an approximate solution  $M_t$  of the assignment problem for  $V$  fixed by matrix balancing (a generalization of Sinkhorn's algorithm).
2. Compute the gradient of the objective with respect to  $V$  by backpropagating through the computations and perform a gradient step.

The overall algorithm, called *XSDC* for "X-Supervised Discriminative Clustering", where "X" can be "un", "semi" or "-" (hence covering all cases), is as follows:

---

Algorithm 1: XSDC

- 1: **Input:** Labeled data  $X_S, Y_S$
- 2: Unlabeled data  $X_U$
- 3: Number of iterations  $T$
- 4: **Initialize:**  $V_1 \leftarrow$  Optimize (1) over  $V$  using  $X_S, Y_S$
- 5: **for**  $t = 1, \dots, T$  **do**
- 6:  $X_t, Y_t \leftarrow$  Draw mini-batch of samples
- 7:  $M_t \leftarrow$  MatrixBalancing( $A(\Phi(X_t; V_t)), Y_t Y_t^T$ )
- 8:  $V_{t+1} \leftarrow$  GradientStep( $\Phi(X_t; V_t), M_t, V_t$ )
- 9: **end for**
- 10:  $\hat{Y}_U \leftarrow$  NearestNeighbor( $\Phi(X; V_T), Y_S$ )
- 11:  $\hat{W}, \hat{b} \leftarrow$  RegLeastSquares( $X; [Y_S, \hat{Y}_U]$ )
- 12: **Output:**  $\hat{Y}_U, V_T, \hat{W}, \hat{b}$

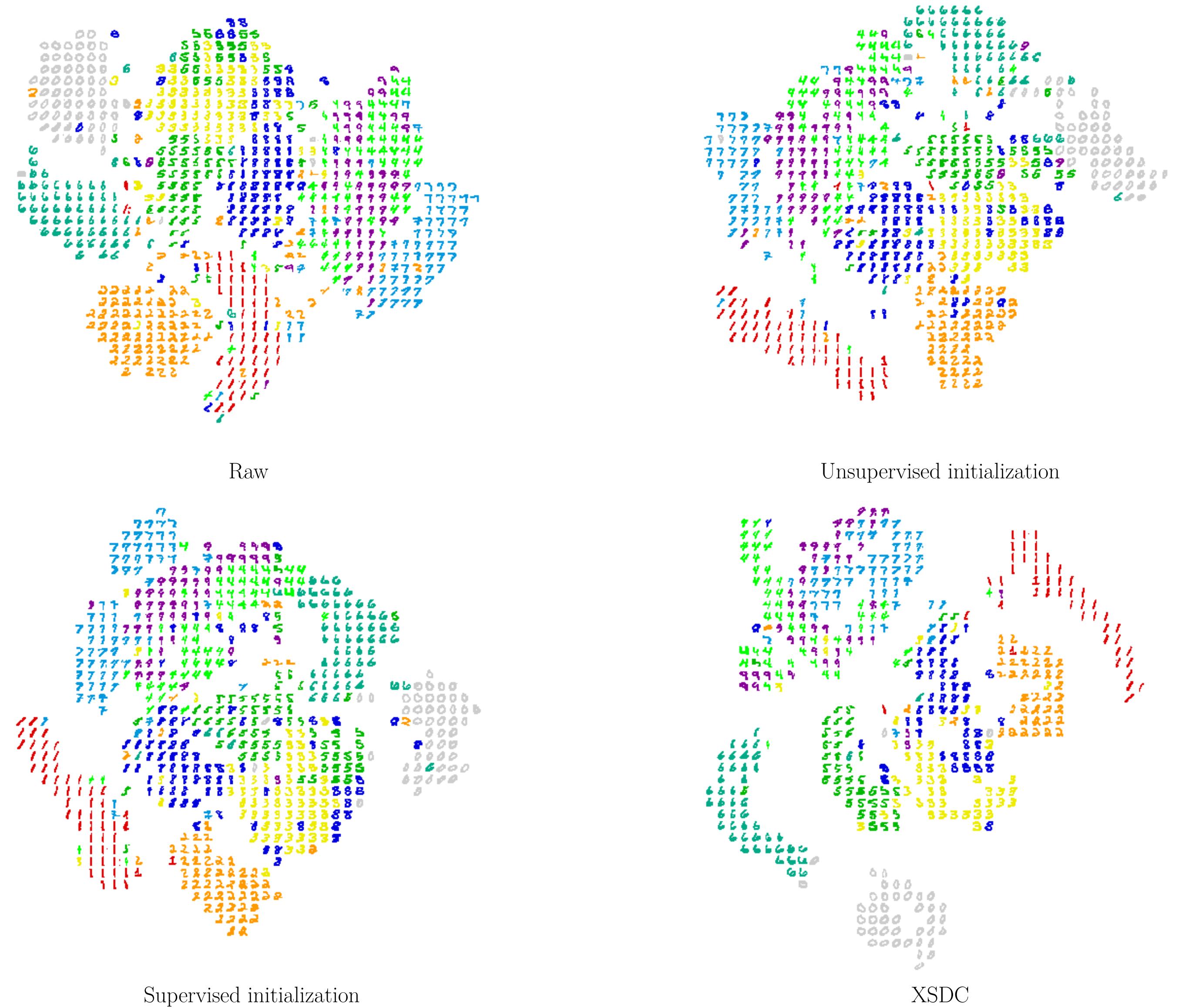
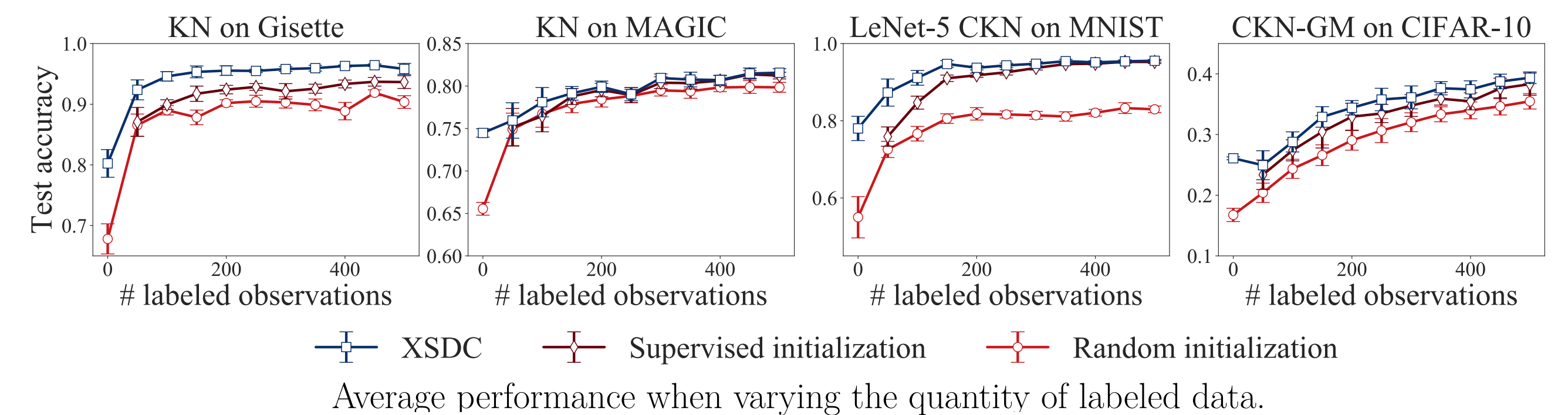
---

## Acknowledgements

This project was supported by the National Science Foundation and the "Learning in Machines and Brains" program of the Canadian Institute for Advanced Research.

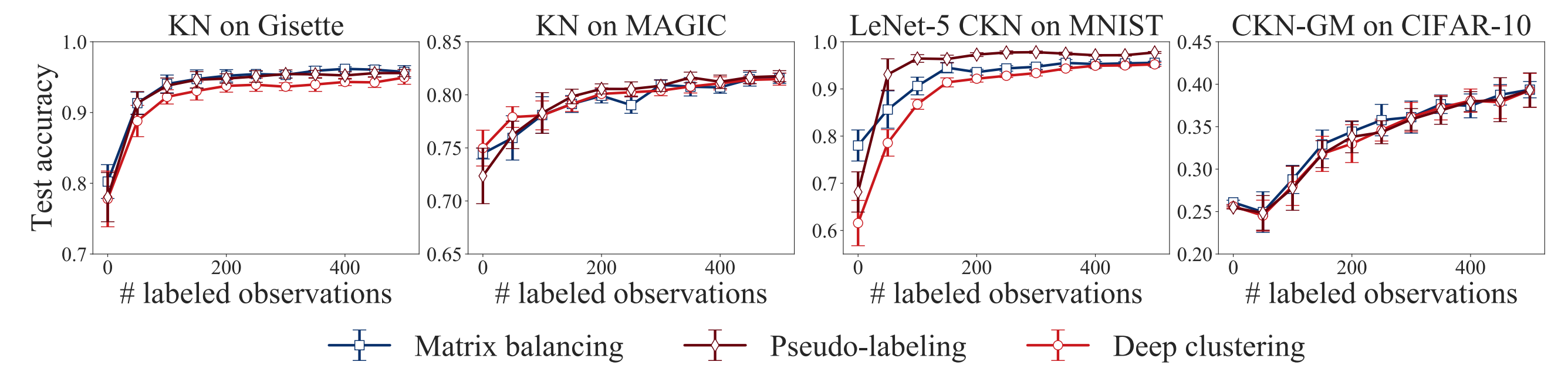
## Results

1. XSDC improves upon using only labeled data in the cases where additional labeled data would help but is unavailable.

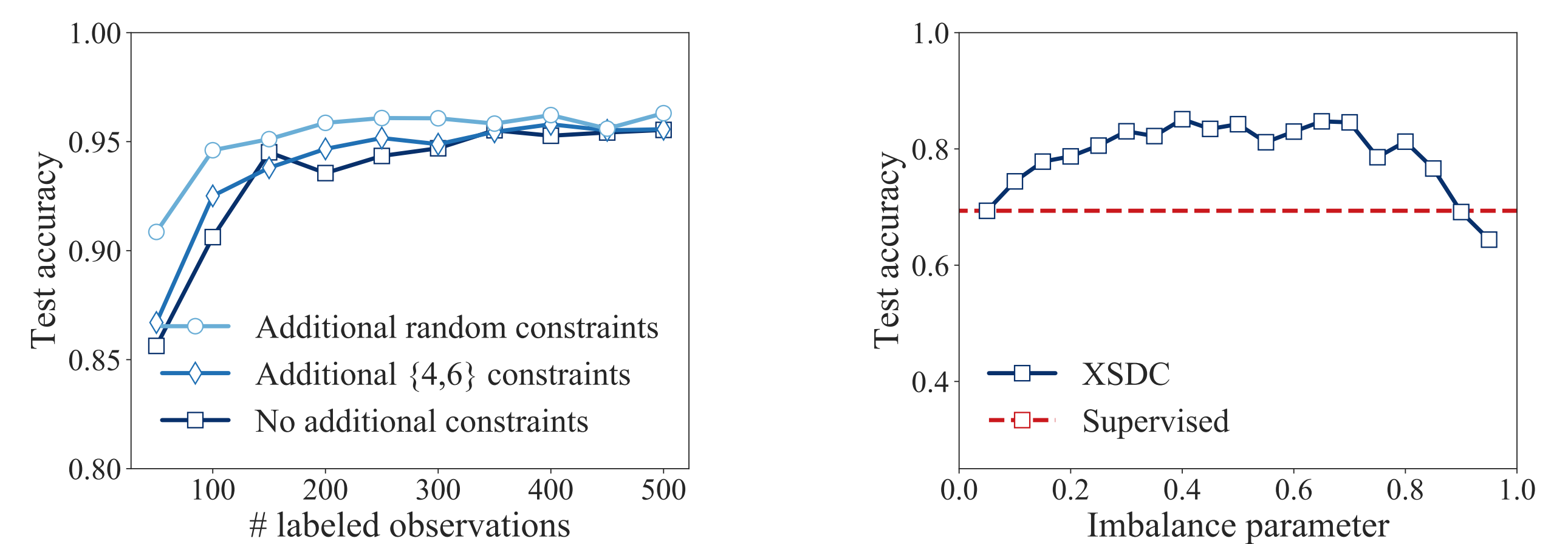


Visualizations of unlabeled MNIST features obtained when training a LeNet-5 CKN with 50 labeled points.

2. The labeling algorithm in XSDC ("matrix balancing") is typically competitive with alternative, less principled, labeling approaches.



3. XSDC can use additional must-link and must-not-link constraints and can handle moderately unbalanced datasets.



## References

- F. R. Bach and Z. Harchaoui. DIFFRAC: a discriminative and flexible framework for clustering. In *Neural Information Processing Systems*, 2007.
- M. Caron, P. Bojanowski, A. Joulin, and M. Douze. Deep clustering for unsupervised learning of visual features. In *European Conference on Computer Vision*, 2018.
- D.-H. Lee. Pseudo-label: The simple and efficient semi-supervised learning method for deep neural networks. In *International Conference on Machine Learning Workshop on Challenges in Representation Learning*, 2013.