

Latent Feature Lasso

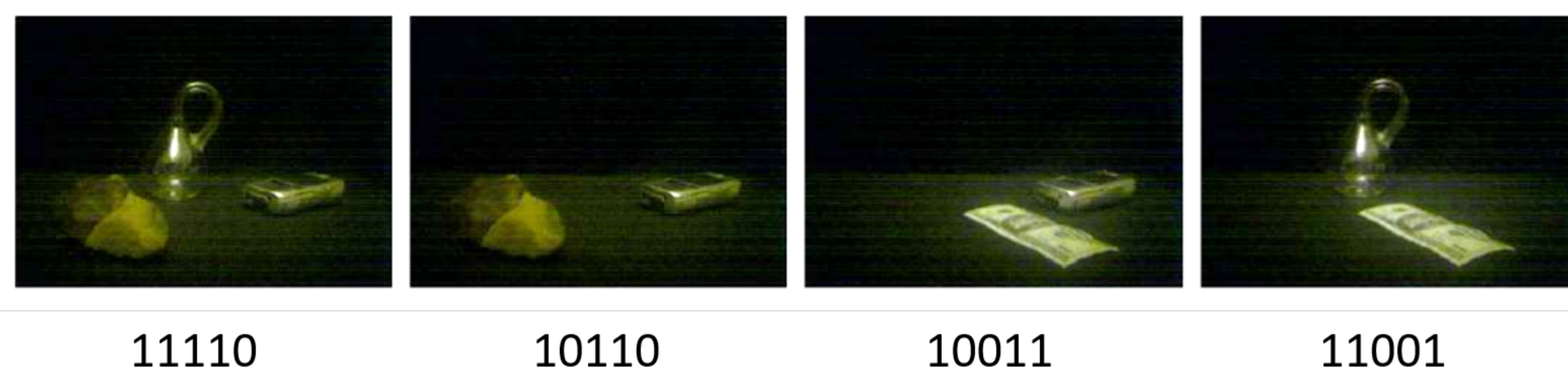
Ian E.H. Yen¹, Wei-Cheng Lee², Sung-En Chang², Arun S. Suggala¹, Shou-De Lin² and Pradeep Ravikumar¹

¹Carnegie Mellon University. ²National Taiwan University

Abstract

- In this work, we propose a novel convex estimator (Latent Feature Lasso) for Latent Feature Model.
- To best of our knowledge, this is the first method with low-order polynomial runtime and sample complexity without restrictive assumptions on the data distribution for LFM.
- In experiments, the Latent Feature Lasso significantly outperforms other methods when there is a larger number of latent features.
- The method enjoys a runtime of $O(ND + DK^2)$ runtime per iter, more scalable than a typical $O(NDK^2)$ of existing approaches.

Latent Feature Models



- Latent Feature Model (LFM) is a generalization of Mixture Model, where each observation is an additive combination of latent features.

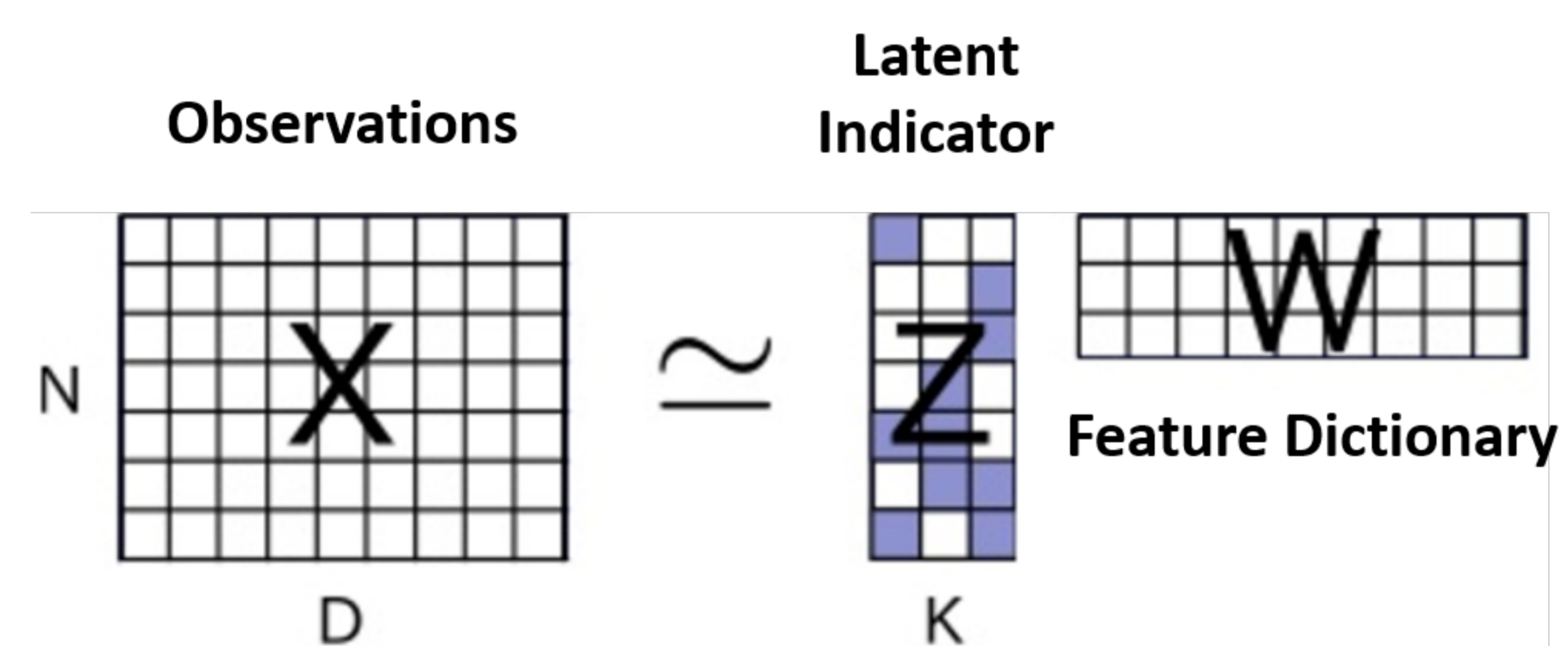
Discriminative	Multiclass Classification	Multilabel Classification
Generative	Mixture Model	Latent Feature Model

- In Latent Feature Model, each observation

$$\mathbf{x}_n = W^T \mathbf{z}_n + \epsilon_n$$

where $\mathbf{x}_n \in \mathbb{R}^D$: observation, $W \in \mathbb{R}^{K \times D}$: feature dictionary, $\mathbf{z}_n \in \{0, 1\}^K$: binary latent indicators, and $\epsilon_n \in \mathbb{R}^D$: noise.

- Mixture Model is a special case with $\|\mathbf{z}_n\|_0 = 1$.



Related Works & Results

- Goal:** Find dictionary $W_{K \times D}$ and latent indicators $Z: N \times K$ that best approximates observation $X: N \times D$.
- Existing Approaches:**
 - MCMC, Variational (Indian Buffet Process): No finite-time guarantee.
 - Spectral Method (Tung 2014): $O(DK^6)$ sample complexity. ($z \sim \text{Ber}(\pi)$, $x \sim N(W^T z, \sigma)$).
 - Matrix Factorization (Slawski et al., 2013): $O(NK2^K)$ runtime complexity for exact recovery (noiseless).
- This Paper:**
 - A convex estimator — Latent Feature Lasso.
 - Low-order polynomial runtime and sample complexity.
 - No restrictive assumption on $p(X)$, even allows model mis-specification.

Convex Formulation via Atomic Norm

- Empirical Risk Minimization:

$$\min_{Z \in \{0,1\}^{N \times K}} \left\{ \min_{W \in \mathbb{R}^{K \times D}} \frac{1}{2N} \|X - ZW\|_F^2 + \frac{\tau}{2} \|W\|_F^2 \right\}$$

- Given Z , the dual problem w.r.t. W is:

$$\min_{M=ZZ^T \in \{0,1\}^{N \times N}} \left\{ \max_{A \in \mathbb{R}^{N \times D}} \frac{-1}{2N^2 \tau} \text{tr}(AA^T M) - \frac{1}{N} \sum_{i=1}^N L^*(x_i, -A_{i,:}) \right\}$$

- Key insight:** the function is convex w.r.t. M .
- Enforce structure $M = ZZ^T$ via an atomic norm.

- Let $\mathcal{S} := \{k \mid \mathbf{z}_k \in \{0, 1\}^N\}$. We define Atomic Norm:

$$\|M\|_{\mathcal{S}} := \min_{c \geq 0} \sum_{k \in \mathcal{S}} c_k \quad \text{s.t.} \quad M = \sum_{k \in \mathcal{S}} c_k \mathbf{z}_k \mathbf{z}_k^T$$

- The Latent Feature Lasso estimator:

$$\min_M g(M) + \lambda \|M\|_{\mathcal{S}}$$

- Equivalently, one can solve the estimator by

$$\min_{c \in \mathbb{R}_+^{|\mathcal{S}|}} g\left(\sum_{k \in \mathcal{S}} c_k \mathbf{z}_k \mathbf{z}_k^T\right) + \lambda \|c\|_1$$

Question: How to optimize with $|\mathcal{S}| = 2^N$ variables?

Greedy Coordinate Descent via MAX-CUT

- At each iteration, we find the coordinate of steepest descent:

$$j^* = \underset{j}{\text{argmax}} -\nabla_j f(c) = \underset{z \in \{0,1\}^N}{\text{argmax}} \langle -\nabla g(M), zz^T \rangle \quad (1)$$

which is a Boolean Quadratic problem similar to MAX-CUT:

$$\max_{z \in \{0,1\}^N} z^T C z$$

- Can be solved to a 3/5-approximation by rounding from a special type of SDP with $O(ND)$ iterative solver.

Active-Set Algorithm

0. $\mathcal{A} = \emptyset$, $c = 0$.

for $t = 1 \dots T$ do

- Find an approximate greedy atom zz^T by MAX-CUT-like problem:

$$\max_{z \in \{0,1\}^N} \langle -\nabla g(M), zz^T \rangle$$

- Add zz^T to an active set \mathcal{A} .

- Refine $c_{\mathcal{A}}$ via Proximal Gradient Method on:

$$\min_{c \geq 0} g\left(\sum_{k \in \mathcal{A}} c_k \mathbf{z}_k \mathbf{z}_k^T\right) + \lambda \|c\|_1$$

- Eliminate $\{\mathbf{z}_k \mathbf{z}_k^T \mid c_k = 0\}$ from \mathcal{A} .

end for.

- Finding approximate greedy coordinate costs $O(ND)$ (via SDP).
- Evaluating $\nabla g(M)$: a least-square problem of cost $O(DK^2)$.
- Each iteration costs $\underbrace{O(ND)}_{\text{MAX-CUT}} + \underbrace{O(DK^2)}_{\text{Least-Square}}$

Runtime Complexity

MCMC	Variational	MF-Binary	BP-Means	Spectral	LatentLasso
$(NDK^2)T$	$(NDK^2)T$	$(NK)2^K$	$(NDK^3)T$	$ND + K^5 \log(K)$	$(ND + K^2 D)T$

Theoretical Results: Risk Bound

Let the population risk of a dictionary W be

$$r(W) := E \left[\min_{z \in \{0,1\}^K} \frac{1}{2} \|x - W^T z\|_2^2 \right]$$

Let W^* be an optimal dictionary of size K , the algorithm outputs \hat{W} with

$$r(\hat{W}) \leq r(W^*) + \epsilon$$

as long as

$$t = \Omega\left(\frac{K}{\epsilon}\right) \quad \text{and} \quad N = \Omega\left(\frac{DK}{\epsilon^3} \log\left(\frac{RK}{\epsilon \rho}\right)\right)$$

- The result trades between risk and sparsity.
- No assumption on x except that of boundedness.
- The sample complexity is (quasi) linear to D and K .

Identifiability

Let $\text{rank}(\Theta^*) = K$. The decomposition $ZW = \Theta^*$ is unique if

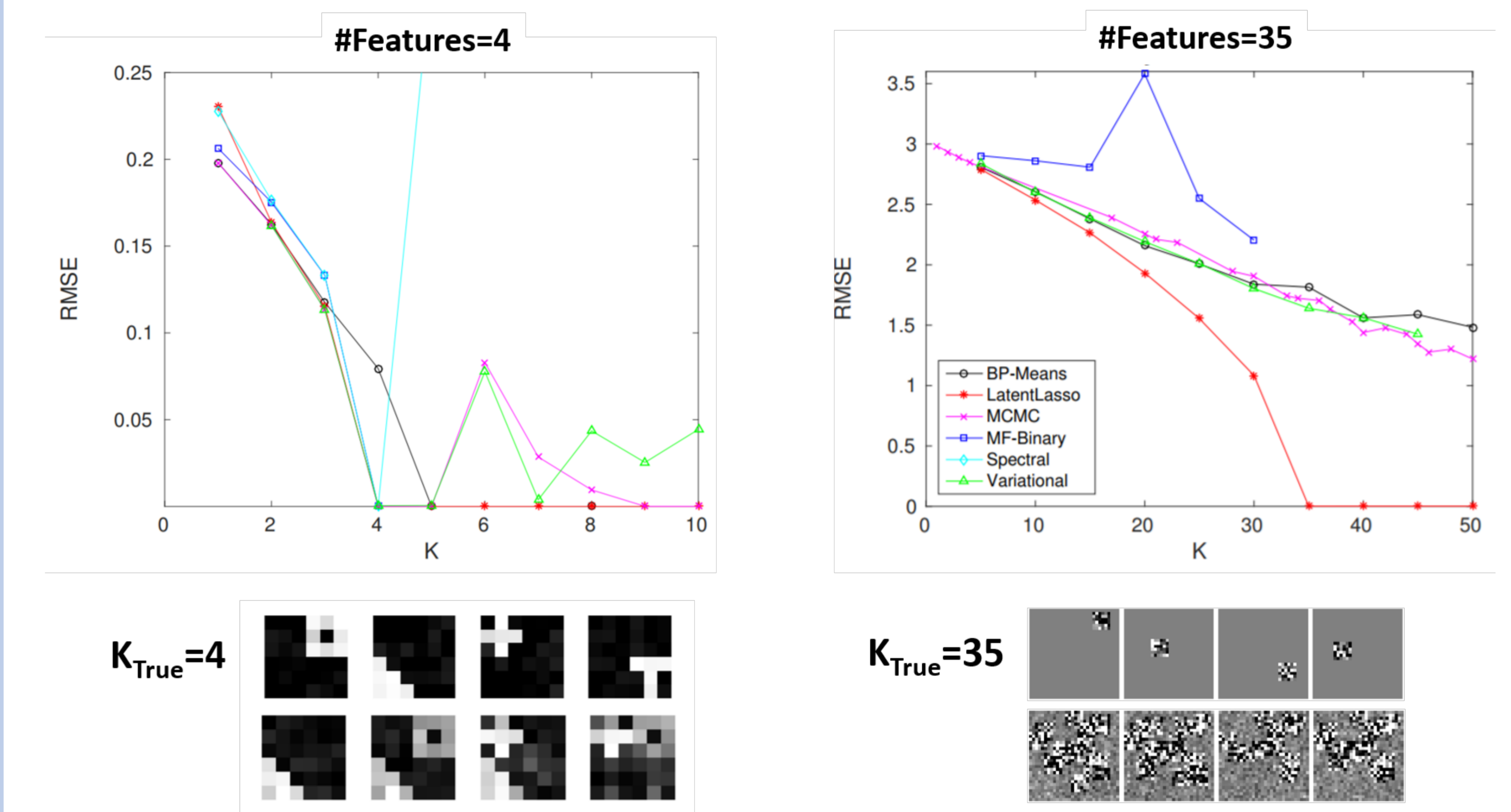
- $Z^*: N \times K$ and $W^*: K \times D$ are both of rank K .
- $\text{span}(Z^*) \cap \{0, 1\}^N \setminus \{0\} = \{Z_{:,j}^*\}_{j=1}^K$.

Theoretical Results: Exact Recovery (noiseless)

Let $X = Z^* W^*$, and (Z_A, W_A) be a solution of Latent Feature Lasso. If the identifiability holds and W_A has full row-rank:

$$\{Z_{:,j}^*\}_{j \in \mathcal{A}} = \{Z_{:,j}^*\}_{j=1}^K, \quad \{W_{j,:}^*\}_{j \in \mathcal{A}} = \{W_{j,:}^*\}_{j=1}^K$$

Experiments on Synthetic Data



Experiments on Real Data

