

Latent Feature Lasso

Ian E.H. Yen^{*}, Wei-Cheng Lee[†], Sung-En Chang[†], Arun Suggala^{*},
Shou-De Lin[†] and Pradeep Ravikumar^{*}.

^{*} Carnegie Mellon University

[†] National Taiwan University

Latent Feature Models



11110



10110



10011



11001

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

Discriminative	Multiclass Classification	Multilabel Classification
Generative	Mixture Model	Latent Feature Model

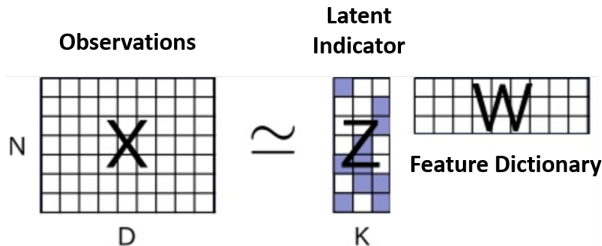
Latent Feature Models

- 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 $\|\mathbf{z}_n\|_0 = 1$.



Latent Feature Models: Result Summary

- **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**.

- 1 Latent-Feature Models
- 2 Latent Feature Lasso—A Convex Estimator
 - Convex Formulation via Atomic Norm
 - Greedy Coordinate Descent via MAX-CUT
- 3 Theoretical Results
- 4 Empirical Results

Latent Feature Model: Estimation

- 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:

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

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

Latent Feature Model: Estimation

- Let $\mathcal{S} := \{\mathbf{z}\mathbf{z}^T \mid \mathbf{z} \in \{0, 1\}^N\}$.
- The "Latent-Feature" Atomic Norm:

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

- The Latent Feature Lasso estimator:

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

- Equivalently, one can solve the estimator by

$$\min_{\mathbf{c} \in \mathbb{R}_+^{|\mathcal{S}|}} g\left(\sum_{k=1}^{2^N} c_k \mathbf{z}_k \mathbf{z}_k^T\right) + \lambda \|\mathbf{c}\|_1$$

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

- 1 Latent-Feature Models
- 2 Latent Feature Lasso—A Convex Estimator
 - Convex Formulation via Atomic Norm
 - Greedy Coordinate Descent via MAX-CUT
- 3 Theoretical Results
- 4 Empirical Results

Greedy Coordinate Descent via MAX-CUT

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

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

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**.

Greedy Coordinate Descent via MAX-CUT

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

for $t = 1 \dots T$ **do**

1. Find an **approximate greedy** atom $\mathbf{z}\mathbf{z}^T$ by MAX-CUT-like problem:

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

2. Add $\mathbf{z}\mathbf{z}^T$ to an **active set** \mathcal{A} .
3. **Refine** $\mathbf{c}_{\mathcal{A}}$ via Proximal Gradient Method on:

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

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

end for.

- Evaluating $\nabla g(M)$ requires solving a **least-square problem** of cost $O(DK^2)$.
- Each iteration costs $\underbrace{O(ND)}_{\text{MAX-CUT}} + \underbrace{O(DK^2)}_{\text{Least-Square}}$

- 1 Latent-Feature Models
- 2 Latent Feature Lasso—A Convex Estimator
 - Convex Formulation via Atomic Norm
 - Greedy Coordinate Descent via MAX-CUT
- 3 Theoretical Results
- 4 Empirical Results

Risk Analysis

Let the **population risk** of a dictionary W be

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

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

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

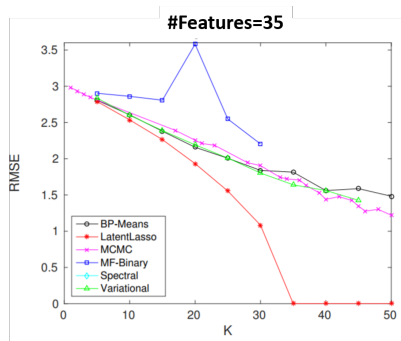
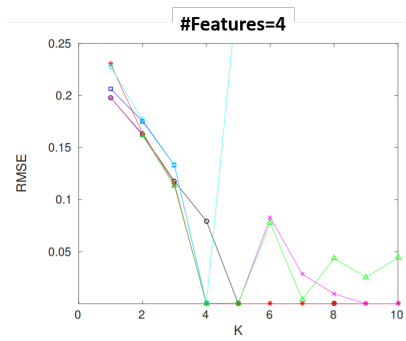
as long as

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

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

- 1 Latent-Feature Models
- 2 Latent Feature Lasso—A Convex Estimator
 - Convex Formulation via Atomic Norm
 - Greedy Coordinate Descent via MAX-CUT
- 3 Theoretical Results
- 4 Empirical Results

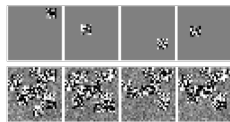
Results on Synthetic Data



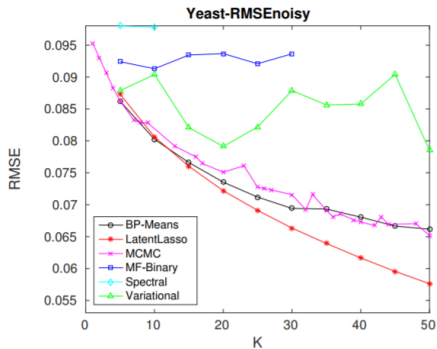
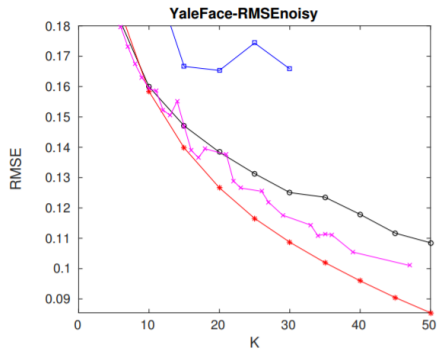
$K_{\text{True}}=4$



$K_{\text{True}}=35$



Results on Real Data



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

- MCMC, Variational, BP-Means take up to 1000s training time, while LatentLasso takes up to 100s.

Conclusion

- In this work, we propose a novel **convex estimator (Latent Feature Lasso)** for the estimation of 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.
- In experiments, the Latent Feature Lasso **significantly outperforms** other methods in terms of **accuracy** and **time**, when there is a **larger number of latent features**.