



Proximal alternating algorithms in dictionary learning

Bao Chenglong

Department of Mathematics

National University of Singapore

Joint work with Yuhui Quan, Hui Ji and Zuowei Shen.

Which kind of data we face?

Authorities are <i>only too aware</i> that Kashgar is 4,000 kilometres (2,500 miles) from Beijing but <i>only</i> a tenth of the distance from the Pakistani border, and are <i>desperate to ensure instability or militancy</i> does not leak over the frontiers.
Taiwan-made products <i>stood a good chance</i> of becoming <i>even more competitive thanks to</i> wider access to overseas markets and lower costs for material imports, he said.
"March <i>appears</i> to be a <i>more reasonable</i> estimate while earlier admission <i>cannot be entirely ruled out</i> ," according to Chen, also Taiwan's chief WTO negotiator.
friday evening plans were great, but saturday's plans <i>didn't go as expected</i> – i went dancing & it was an <i>ok</i> club, but <i>terribly crowded</i> :-/
WHY THE <i>HELL</i> DO YOU GUYS ALL HAVE MRS. KENNEDY! SHES A FUCKING DOUCHE
AT&T was <i>okay</i> but whenever they do something <i>nice</i> in the name of customer service it seems like a favor, while T-Mobile makes that a <i>normal everyday thing</i>
obama should be <i>impeached</i> on <i>TREASON</i> charges. Our Nuclear arsenal was TOP Secret. Till HE told our enemies what we had. <i>#Coward #Traitor</i>
My graduation speech: "I'd like to <i>thanks</i> Google, Wikipedia and my computer! <i>:D #iThingteens</i>

Text data

Which kind of data we face?

Authorities are <i>only too aware</i> that Kashgar is 4,000 kilometres (2,500 miles) from Beijing but <i>only</i> a tenth of the distance from the Pakistani border, and are <i>desperate</i> to <i>ensure instability or militancy</i> does not leak over the frontiers.
Taiwan-made products <i>stood a good chance</i> of becoming <i>even more competitive thanks to</i> wider access to overseas markets and lower costs for material imports, he said.
"March <i>appears</i> to be a <i>more reasonable</i> estimate while earlier admission <i>cannot be entirely ruled out</i> ," according to Chen, also Taiwan's chief WTO negotiator.
friday evening plans were great, but saturday's plans <i>didn't go as expected</i> – i went dancing & it was an <i>ok</i> club, but <i>terribly crowded</i> :-/
WHY THE HELL DO YOU GUYS ALL HAVE MRS. KENNEDY! SHES A FUCKING DOUCHE
AT&T was <i>okay</i> but whenever they do something <i>nice</i> in the name of customer service it seems like a favor, while T-Mobile makes that a <i>normal everyday thin</i>
obama should be <i>impeached</i> on TREASON charges. Our Nuclear arsenal was TOP Secret. Till HE told our enemies what we had. #Coward#Traitor
My graduation speech: "I'd like to <i>thanks</i> Google, Wikipedia and my computer! :D #iThingteens

Text data



Healthy data



Financial data

Which kind of data we face?

Authorities are <i>only too aware</i> that Kashgar is 4,000 kilometres (2,500 miles) from Beijing but <i>only</i> a tenth of the distance from the Pakistani border, and are <i>desperate to ensure instability or militancy</i> does not leak over the frontiers.
Taiwan-made products <i>stood a good chance</i> of becoming <i>even more competitive thanks to</i> wider access to overseas markets and lower costs for material imports, he said.
"March <i>appears</i> to be a <i>more reasonable</i> estimate while earlier admission <i>cannot be entirely ruled out</i> ," according to Chen, also Taiwan's chief WTO negotiator.
friday evening plans were great, but saturday's plans <i>didn't go as expected</i> – i went dancing & it was an <i>ok</i> club, but <i>terribly crowded</i> :-/
WHY THE HELL DO YOU GUYS ALL HAVE MRS. KENNEDY! SHES A FUCKING DOUCHE
AT&T was <i>okay</i> but whenever they do something <i>nice</i> in the name of customer service it seems like a favor, while T-Mobile makes that a <i>normal everyday thin</i>
obama should be <i>impeached</i> on <i>TREASON</i> charges. Our Nuclear arsenal was TOP Secret. Till HE told our enemies what we had. <i>#Coward#Traitor</i>
My graduation speech: "I'd like to <i>thanks</i> Google, Wikipedia and my computer! :D #iThingteens

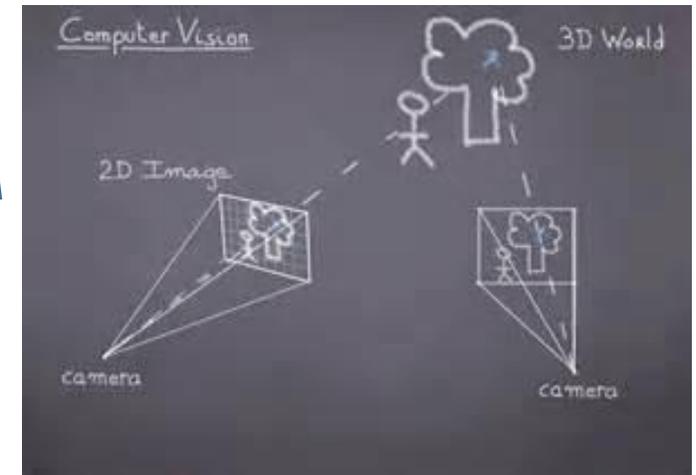
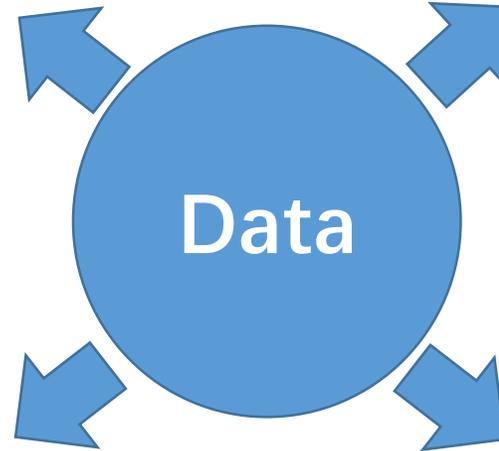
Text data



Healthy data

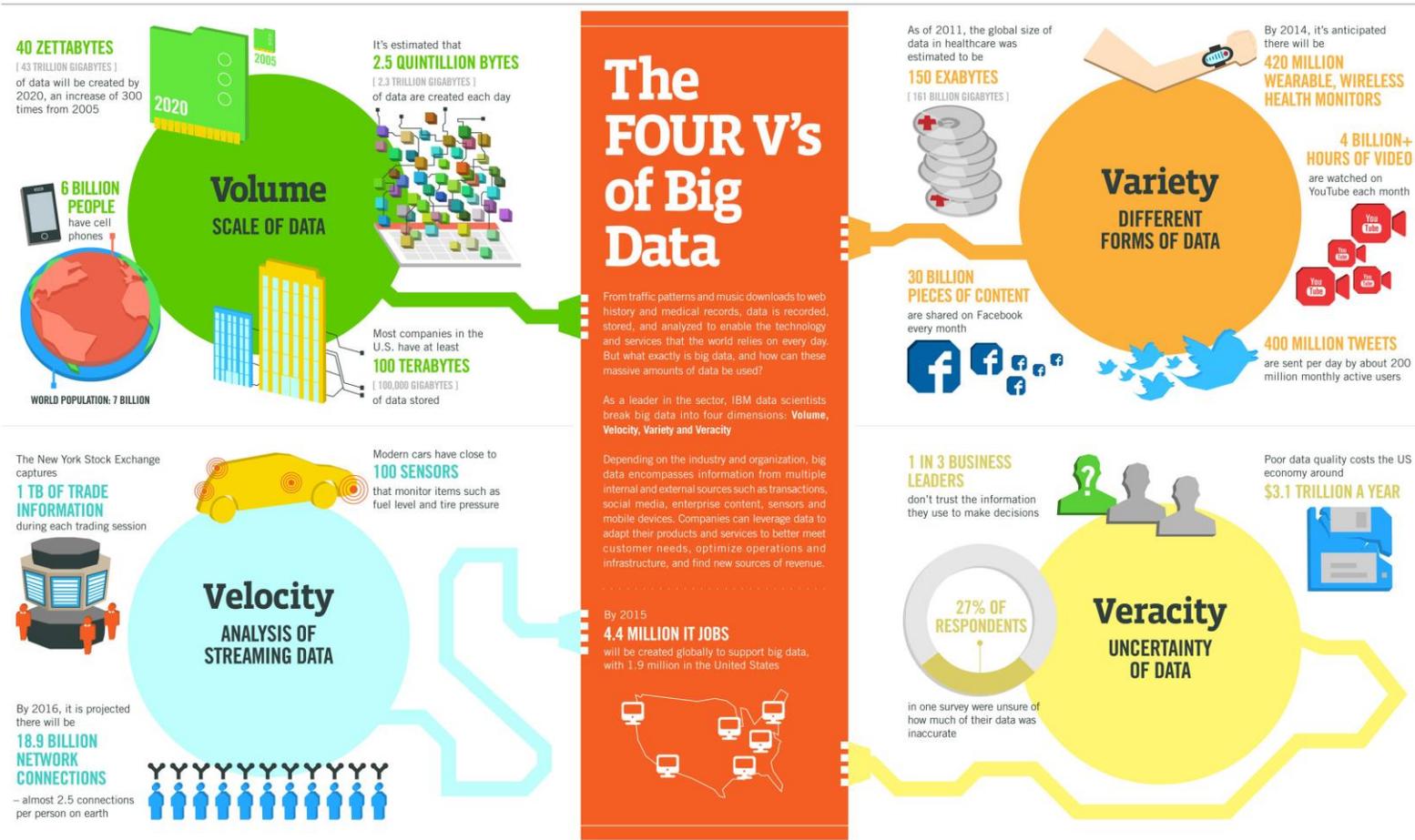


Financial data



Vision data

Why focus on data representation?



- 4Vs in data science
 - a) Volume
 - b) Variety
 - c) Velocity
 - d) Veracity

- Key of the success
 - Efficient and effective data representation

Data representation

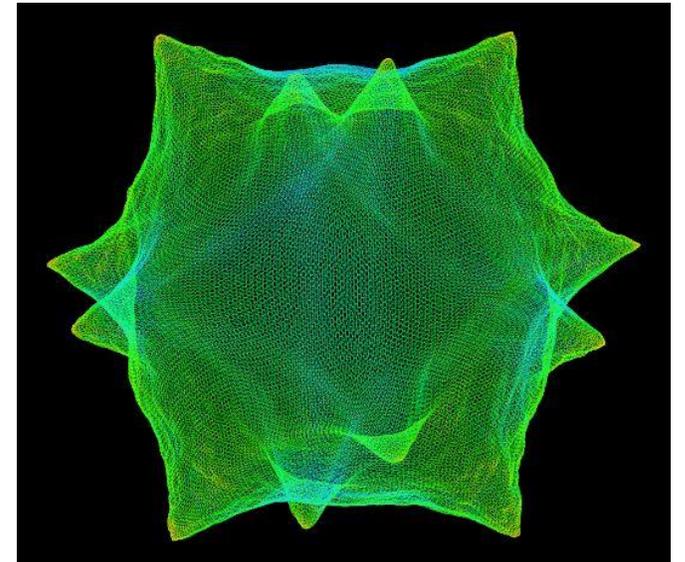
- Main difficulty: curse of dimensionality (e.g. Hughes phenomenon)

- Sparsity prior
 - Regularity
 - Temporal and spatial coherence
 - Hierarchical organization of explanatory factors
 - **Dictionary learning**



Non-convex models

- Goal: “good” sparse representation
 - Efficient and convergent numerical schemes and effective models

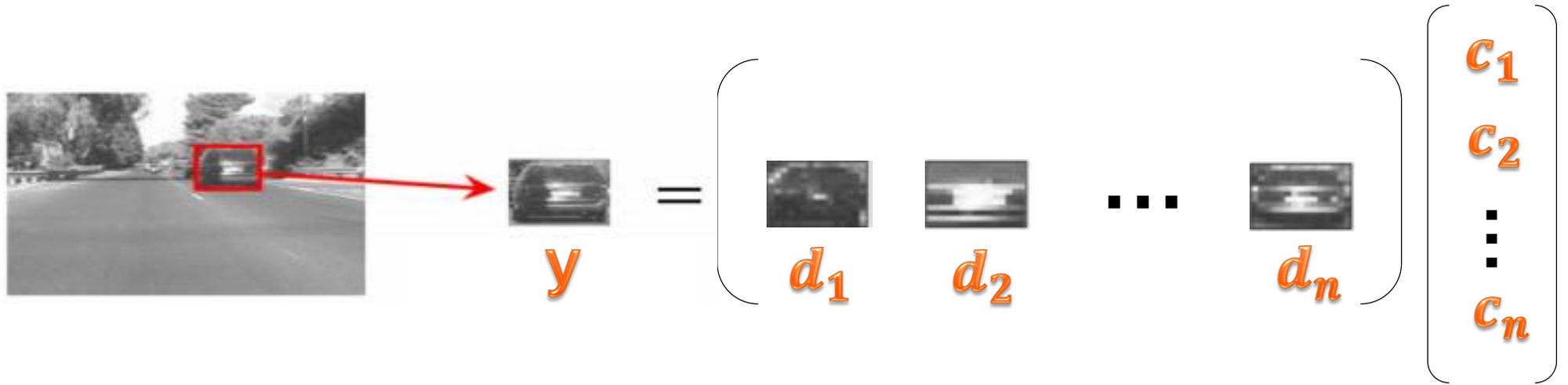


Sparse matrices

Sparsity based dictionary learning

- Initially proposed in [Olshausen et al, Nature, 1996]
- Many variants
 - Redundant dictionary learning [Elad et al, TIP, 2006]
 - Data-driven tight frame construction [Cai et al, ACHA, 2014]
 - Task-driven dictionary learning [Mairal et al, TPAMI, 2012]
 - Hierarchical dictionary learning [Jenatton et al, JMLR, 2011]
- Success in applications
 - Image denoising [Vincent et al, JMLR, 2010]
 - Abnormal event detection [Zhao et al, CVPR 2011]
 - Natural language processing [Bagnell et al, NIPS, 2009]
 - Visual tracking [Mei et al, TPAMI, 2011]

Dictionary learning

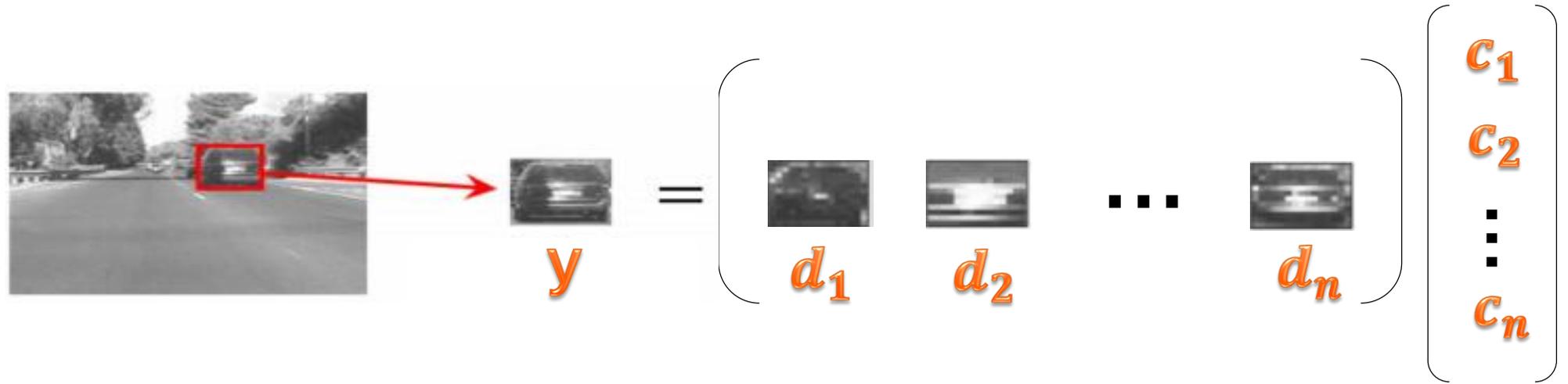


- Data-driven sparse representation
 - Dictionary learning: find an adaptive dictionary $D = (d_1, d_2, \dots, d_n)$ such that

$$y = \sum_i c_i d_i,$$

and **most c_i s are zero.**

Dictionary learning



- Data-driven sparse representation
 - Dictionary learning: find an adaptive dictionary $D = (d_1, d_2, \dots, d_n)$ such that

and **most c_i s are zero**

$$y = \sum_i c_i d_i,$$

How to solve the non-convex model?

The most popular dictionary learning method

- Data $Y = (y_1, y_2 \dots y_p)$ and $D = (d_1, d_2, \dots, d_m)$ and $C = (c_1, c_2, \dots, c_p)$

- K-SVD [Elad et al, TIP, 2006]

$$\min_{D,C} \|Y - DC\|_F^2, s.t. \|c_i\|_0 \leq k, \|d_j\|_2 = 1, \forall i \in [p], \forall j \in [m].$$

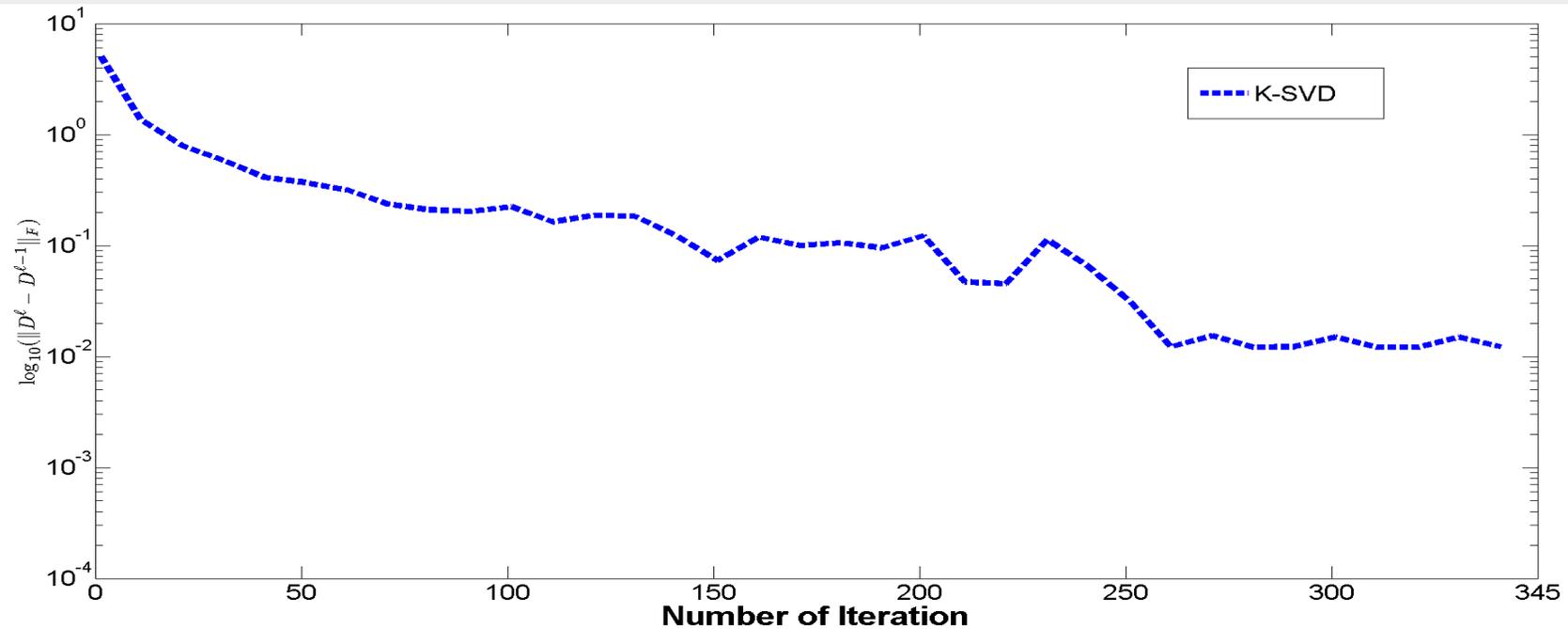
- Alternating minimization: OMP and sequential SVD (No convergence guarantee and high computational complexity)

- K-SVD based applications

- Image classification [Jiang et al, TPAMI, 2013]
- Text corpora representation [Jenatton et al, JMLR, 2011]
- Multi-task learning [Ruvolo et al, ICML, 2014]
- Recommendation system [Gediminas et al, TKDE, 2012]

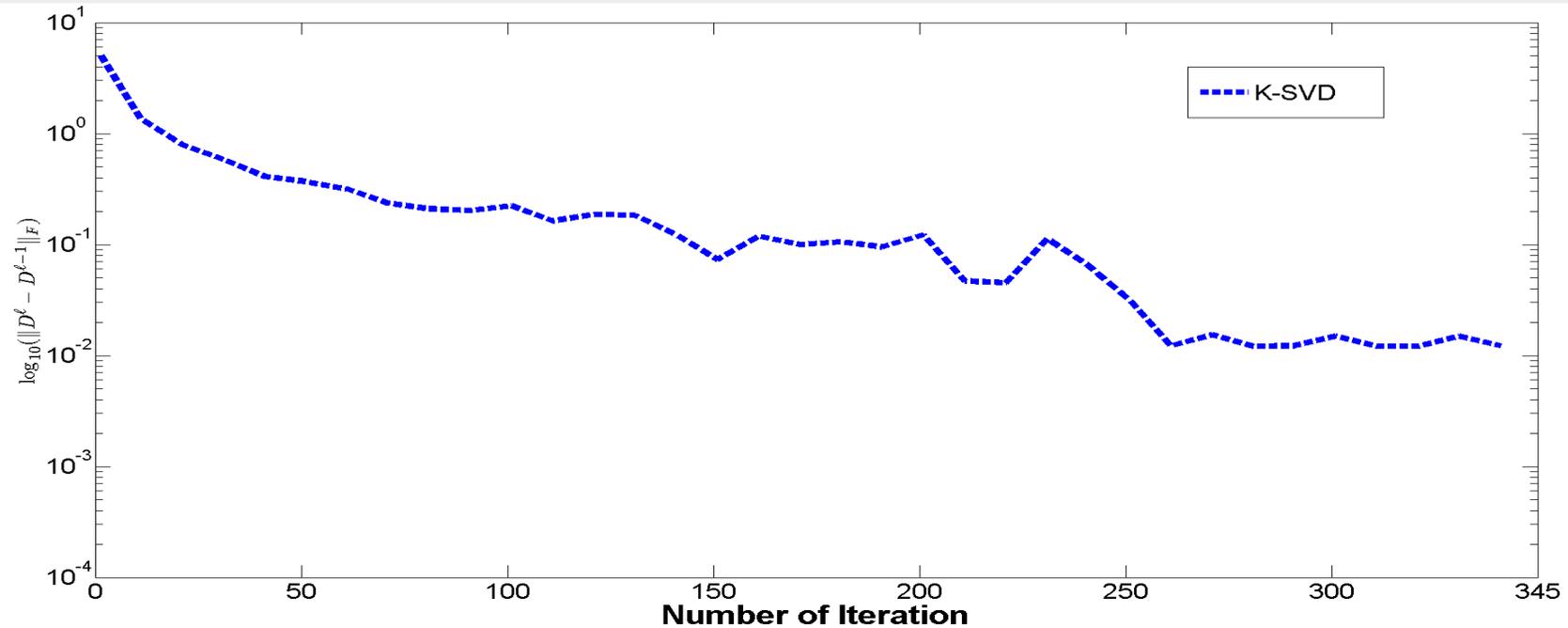
The numerical behavior of K-SVD

L2 norm of increments of the sequence generated by K-SVD



The numerical behavior of K-SVD

L2 norm of increments of the sequence generated by K-SVD



Goal: convergent and efficient numerical algorithm.

Mathematical formulation

- General dictionary learning model

$$\min_{D,C} \sum_{i=1}^p f(y_i, Dc_i) + r_1(C) + r_2(D)$$

- f is the loss function, e.g. ℓ_2 loss,
- r_1, r_2 are constraint functions on C, D , respectively.

Mathematical formulation

- General dictionary learning model

$$\min_{D, C} \sum_{i=1}^p f(y_i, Dc_i) + r_1(C) + r_2(D)$$

- f is the loss function, e.g. ℓ_2 loss,
- r_1, r_2 are constraint functions on C, D , respectively.

- Multi-block non-convex optimization

$$\min_{x=(x_1, \dots, x_n)} H(x) := f(x) + \sum_{i=1}^n r_i(x_i)$$

- Assumptions

1. ∇f is Lipschitz on any bounded set and $\nabla_i f$ is L_i -Lipschitz
2. $r_i, i \in [n]$ are proper and lower semi-continuous

Existing schemes

- Define $f_i^k(\cdot) := f(x_1^{k+1}, \dots, x_{i-1}^{k+1}, \cdot, x_{i+1}^k, \dots, x_n^k)$

Alternating minimization

$$x_i^{k+1} \in \operatorname{argmin} f_i^k(x_i) + r_i(x_i)$$

Proximal alternating minimization

$$x_i^{k+1} \in \operatorname{argmin} f_i^k(x_i) + r_i(x_i) + \lambda_i \|x_i - x_i^k\|_F^2$$

Proximal alternating linearized minimization

$$x_i^{k+1} \in \operatorname{argmin} \langle \nabla f_i^k(x_i^k), x_i - x_i^k \rangle + r_i(x_i) + \lambda_i \|x_i - x_i^k\|_F^2$$

Literature review

- Alternating minimization (**AM**)
 - No convergence guarantee
- Proximal alternating minimization (**PAM**)
 - Global convergence property [Attouch et al, MOR, 2010]
- Proximal alternating linearized minimization (**PALM**)
 - Global convergence property [Bolte et al, MP, 2014]
- Hybrid method
 - Multi-convex case [Xu et al, SIIMS, 2013]

Literature review

- Alternating minimization (**AM**)
 - No convergence guarantee
- Proximal alternating minimization (**PAM**)
 - Global convergence property [Attouch et al, MOR, 2010]
- Proximal alternating linearized minimization (**PALM**)
 - Global convergence property [Bolte et al, MP, 2014]
- Hybrid method
 - Multi-convex case [Xu et al, SIIMS, 2013]

Inner iter. No. : **AM** \approx **PAM** $>$ **PALM**

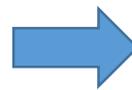
Outer iter. No. : **AM** $<$ **PAM** $<$ **PALM**

Literature review

- Alternating minimization (**AM**)
 - No convergence guarantee
- Proximal alternating minimization (**PAM**)
 - Global convergence property [Attouch et al, MOR, 2010]
- Proximal alternating linearized minimization (**PALM**)
 - Global convergence property [Bolte et al, MP, 2014]
- Hybrid method
 - Multi-convex case [Xu et al, SIIMS, 2013]

Inner iter. No. : AM \approx PAM > PALM

Outer iter. No. : AM < PAM < PALM



Hybrid them for non-convex case.

Hybrid proximal alternating method

- HPAM [Bao et al, TPAMI, 2016]

For $k = 1, 2, \dots$,

for $i = 1, 2, \dots, n$

$$x_i^{k+1} \in \begin{cases} \operatorname{argmin} f_i^k(x_i) + r_i(x_i) + \lambda_i \|x_i - x_i^k\|_F^2 & \text{or} \\ \operatorname{argmin} \langle \nabla f_i^k(x_i^k), x_i - x_i^k \rangle + r_i(x_i) + \lambda_i \|x_i - x_i^k\|_F^2 \end{cases}$$

PAM

PALM

where λ_i is the appropriate step size.

Hybrid proximal alternating method

- HPAM [Bao et al, TPAMI, 2016]

For $k = 1, 2, \dots,$

for $i = 1, 2, \dots, n$

$$x_i^{k+1} \in \begin{cases} \operatorname{argmin} f_i^k(x_i) + r_i(x_i) + \lambda_i \|x_i - x_i^k\|_F^2 & \text{or} \\ \operatorname{argmin} \langle \nabla f_i^k(x_i^k), x_i - x_i^k \rangle + r_i(x_i) + \lambda_i \|x_i - x_i^k\|_F^2 \end{cases}$$

PAM

PALM

where λ_i is the appropriate step size.

Theorem (Global convergence): Let $\{x^k\}$ be the infinite sequence generated by the HPAM. If H is a **KL function** and $\{x^k\}$ is **bounded**. Then, the sequence $\{x^k\}$ converges to a point x^* , which is a critical point of H , i.e. $0 \in \partial H(x^*)$.

Remark: **AM** can be included under certain condition.

Sketch of the proof

- The proof is based on the result from [Attouch et al, MP, 2013].
- Four main steps:

1. Sufficient decrease property

$$H(x^k) - H(x^{k+1}) \geq \rho_1 \|x^k - x^{k+1}\|_F^2, \text{ for some } \rho_1 > 0$$

2. Bounded the subgradient:

$$\text{dist}\left(0, \partial H(x^k)\right) \leq \rho_2 \|x^k - x^{k-1}\|, \text{ for some } \rho_2 > 0$$

3. Subsequence continuity: \exists a subsequence $\{x^{k_j}\}$ such that

$$x^{k_j} \rightarrow \bar{x}, \quad \text{and} \quad H(x^{k_j}) \rightarrow H(\bar{x}).$$

4. H satisfies Kurdyka-Lojasiewicz (KL) property.

HPAM in dictionary learning

- Plain dictionary learning model

$$\min_{D,C} \|Y - DC^T\|_F^2 + \lambda \|C\|_0, \text{ s. t. } \|C\|_\infty \leq M, \|d_j\|_2 = 1, \forall j \in [m].$$

- Numerical schemes

- C1 $(x_1, x_2, \dots, x_n) = (c_1, c_2, \dots, c_q, d_1, d_2, \dots, d_q)$
- C2 $(x_1, x_2, \dots, x_n) = (C, D)$
- C3 $(x_1, x_2, \dots, x_n) = (C, d_1, c_1, d_2, c_2, \dots, d_q, c_q)$

HPAM in dictionary learning

- Plain dictionary learning model

$$\min_{D,C} \|Y - DC^T\|_F^2 + \lambda \|C\|_0, \text{ s. t. } \|C\|_\infty \leq M, \|d_j\|_2 = 1, \forall j \in [m].$$

- Numerical schemes

- C1 $(x_1, x_2, \dots, x_n) = (c_1, c_2, \dots, c_q, d_1, d_2, \dots, d_q)$
- C2 $(x_1, x_2, \dots, x_n) = (C, D)$
- C3 $(x_1, x_2, \dots, x_n) = (C, d_1, c_1, d_2, c_2, \dots, d_q, c_q)$

SCHEMES	S1	S2	S3
Block choice	C1	C2	C3
Step I	PAM	PALM	PALM
Step II	PAM	PALM	PAM+AM

HPAM in dictionary learning

- Plain dictionary learning model

$$\min_{D,C} \|Y - DC^T\|_F^2 + \lambda \|C\|_0, s. t. \|C\|_\infty \leq M, \|d_j\|_2 = 1, \forall j \in [m].$$

- Numerical schemes

- C1 $(x_1, x_2, \dots, x_n) = (c_1, c_2, \dots, c_q, d_1, d_2, \dots, d_q)$
- C2 $(x_1, x_2, \dots, x_n) = (C, D)$
- C3 $(x_1, x_2, \dots, x_n) = (C, d_1, c_1, d_2, c_2, \dots, d_q, c_q)$

SCHEMES	S1	S2	S3
Block choice	C1	C2	C3
Step I	PAM	PALM	PALM
Step II	PAM	PALM	PAM+AM

All the above schemes converge to a critical point.

Scheme 1 (S1)

Block choice: $C1 (x_1, x_2, \dots, x_n) = (c_1, c_2, \dots, c_q, d_1, d_2, \dots, d_q)$

Step 1. Sparse coding: for $i = 1, \dots, m$,

$$c_i^k \in \underset{\|c\|_\infty \leq M}{\operatorname{argmin}} \lambda \|c\|_0 + \|J_i^k - d_i^k c\|_F^2 + \lambda_c \|c - c_i^{k-1}\|_F^2$$

where $J_i^k = Y - \sum_{j < i} d_j^k c_j^{k\top} + \sum_{j > i} d_j^k c_j^{k-1\top}$.

Step 2. Dictionary update: for $i = 1, \dots, m$,

$$d_i^{k+1} \in \underset{d}{\operatorname{argmin}} \|E_i^k - d c_i^{k\top}\|_F^2 + \lambda_d \|d - d_i^k\|_F^2 \text{ s.t. } \|d\| = 1$$

where $E_i^k = Y - \sum_{j < i} d_j^{k+1} c_j^{k\top} + \sum_{j > i} d_j^k c_j^{k\top}$.

Scheme 2 (S2)

Block choice: $C2(x_1, x_2, \dots, x_n) = (C, d_1, d_2, \dots, d_q)$

Step 1. Sparse coding:

$$C^k \in \underset{\|C\|_\infty \leq M}{\operatorname{argmin}} \lambda \|C\|_0 + \langle \nabla_C f(D^k, C^{k-1}), C - C^{k-1} \rangle + \lambda_C \|C - C^{k-1}\|_F^2$$

Step 2. Dictionary update: for $i = 1, \dots, m$,

$$d_i^{k+1} \in \underset{d}{\operatorname{argmin}} \|E_i^k - d c_i^{k\top}\|_F^2 + \lambda_d \|d - d_i^k\|_F^2 \text{ s.t. } \|d\| = 1$$

where $E_i^k = Y - \sum_{j < i} d_j^{k+1} c_j^{k\top} + \sum_{j > i} d_j^k c_j^{k\top}$.

Scheme 3 (S3)

Step 1. Sparse coding:

$$C^k \in \underset{\|C\|_\infty \leq M}{\operatorname{argmin}} \lambda \|C\|_0 + \langle \nabla_C f(D^k, C^{k-1}), C - C^{k-1} \rangle + \lambda_c \|C - C^{k-1}\|_F^2$$

Step 2. Dictionary update: for $i = 1, \dots, m$,

1. Update the d_i

$$d_i^{k+1} \in \operatorname{argmin} \|E_i^k - d c_i^{k\top}\|_F^2 + \lambda_d \|d - d_i^k\|_F^2 \text{ s.t. } \|d\| = 1$$

$$\text{where } E_i^k \stackrel{d}{=} Y - \sum_{j < i} d_j^{k+1} c_j^{k\top} + \sum_{j > i} d_j^k c_j^{k\top}.$$

2. Reupdate the non-zero coefficients c_i

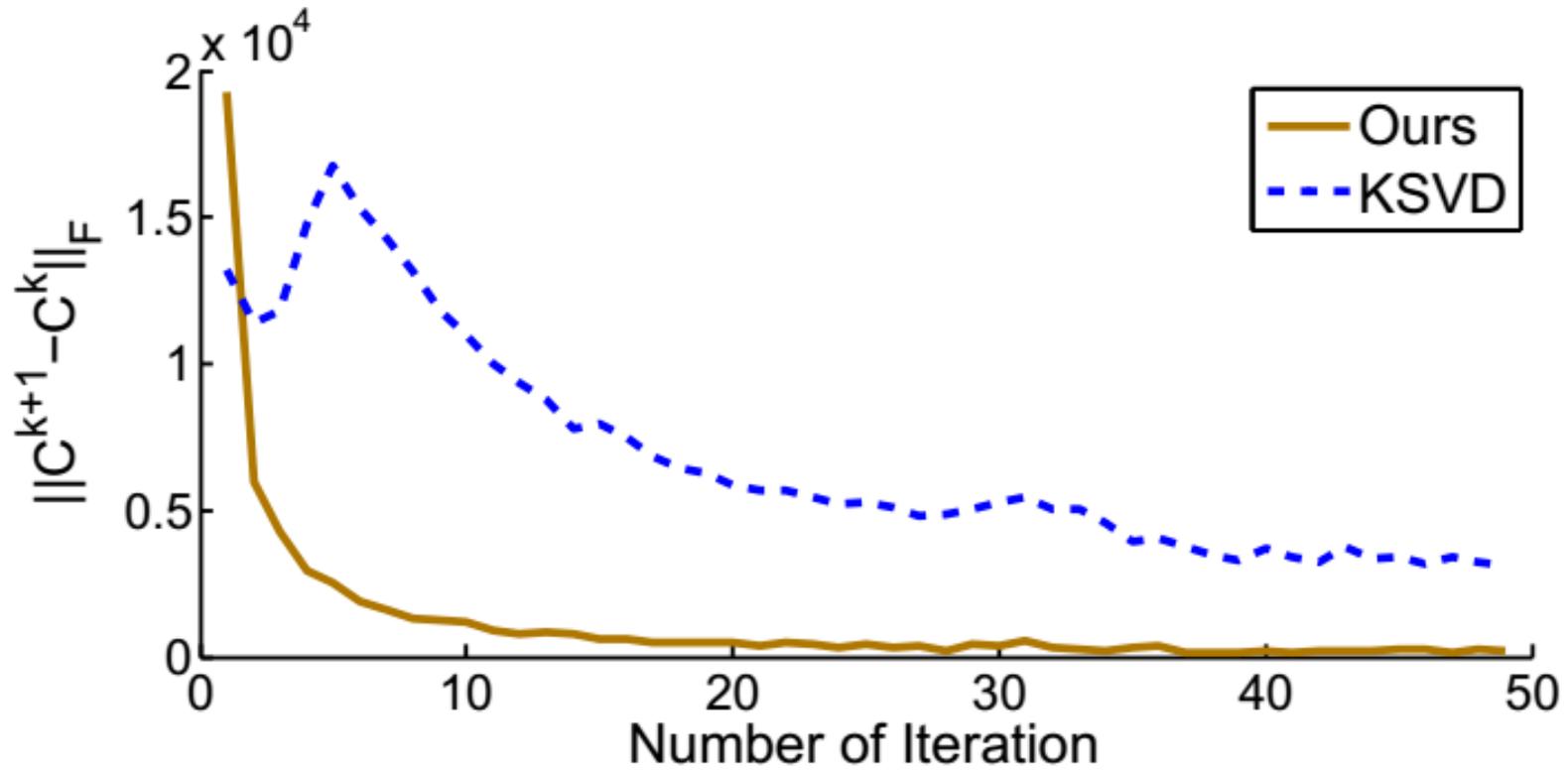
Further decrease the residual.

$$c_i^{k+1} \in \operatorname{argmin} \|E_i^k - d_i^{k+1} c\|_F^2 \text{ s.t. } c(j) = 0, \forall j \in I_j$$

$$\text{where } I_j = \{q: c_i^k(q) = 0\}.$$

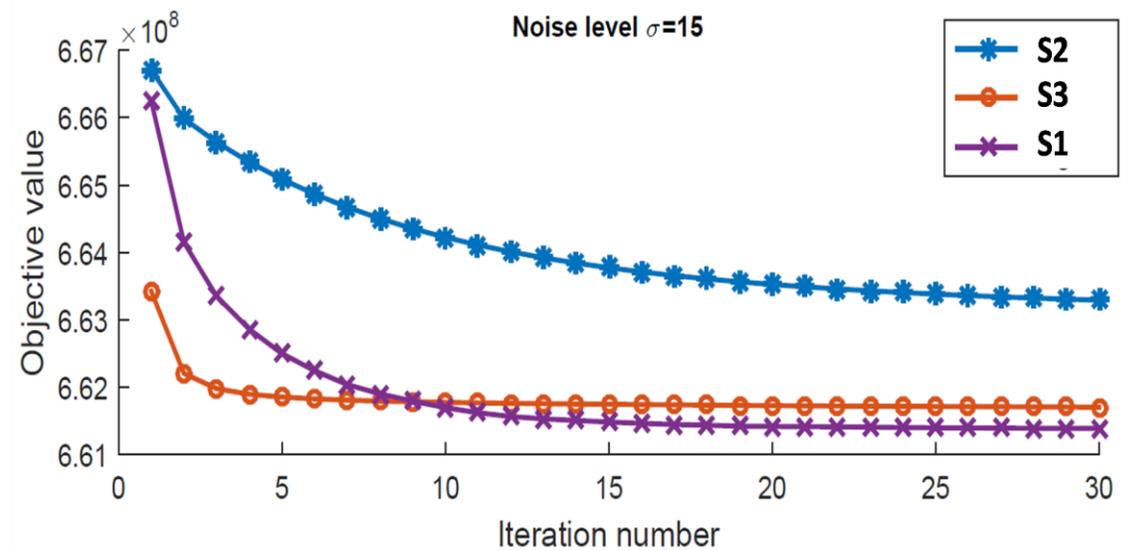
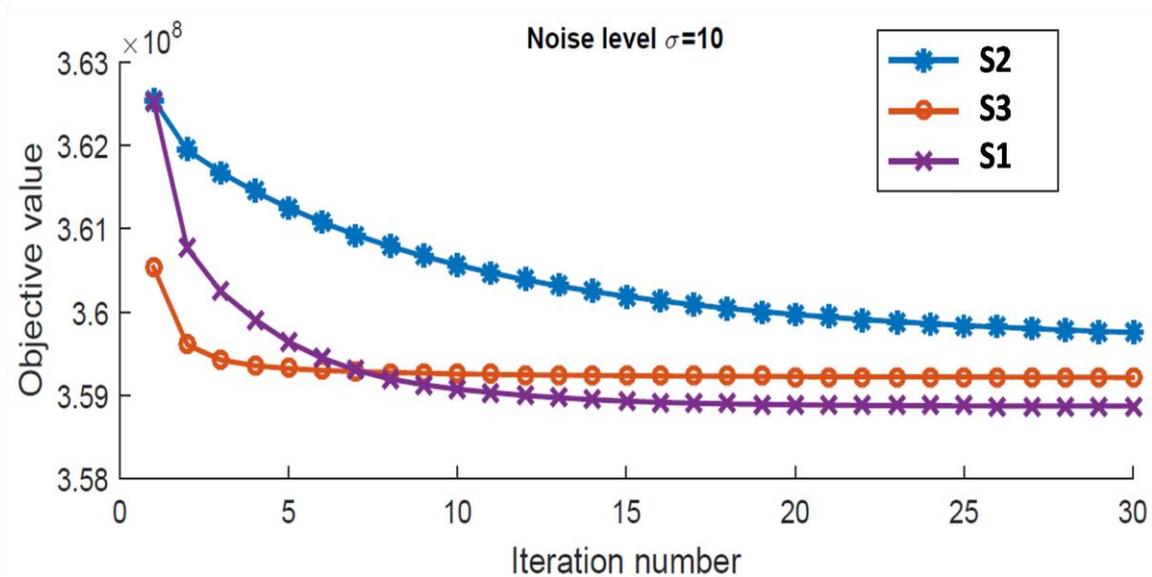
Numerical behavior

Increments of the sequence $\{C^k\}$



Numerical behavior

Objective value versus iteration



Different convergent algorithms achieve different critical points in DL.

Computational efficiency

Training time versus the atom dimensions

Atom Dim.	6*6	8*8	10*10	12*12	14*14	16*16
K-SVD	39	70	114	164	228	308
S1	71	217	465	1011	1848	3094
S2	9	16	28	42	60	86
S3	10	18	30	45	66	96

Efficiency: $S2 \approx S3 > K-SVD > S1$

Image denoising

Fingerprint					Lena				
Image	5	10	15	20	Image	5	10	15	20
Noise	5	10	15	20	Noise	5	10	15	20
K-SVD	36.59	32.39	30.06	28.47	K-SVD	38.59	35.47	33.70	32.38
S1	36.58	32.27	29.87	28.24	S1	38.48	35.37	33.55	32.21
S2	36.50	32.21	29.84	28.18	S2	38.46	35.35	33.50	32.15
S3	36.59	32.35	30.03	28.44	S3	38.49	35.41	33.57	32.25

Performance: $S3 \approx S1 \approx K\text{-SVD} > S2$

Image denoising

Fingerprint					Lena				
Image	5	10	15	20	Image	5	10	15	20
Noise	5	10	15	20	Noise	5	10	15	20
K-SVD	36.59	32.39	30.06	28.47	K-SVD	38.59	35.47	33.70	32.38
S1	36.58	32.27	29.87	28.24	S1	38.48	35.37	33.55	32.21
S2	36.50	32.21	29.84	28.18	S2	38.46	35.35	33.50	32.15
S3	36.59	32.35	30.03	28.44	S3	38.49	35.41	33.57	32.25

Performance: $S3 \approx S1 \approx K-SVD > S2$

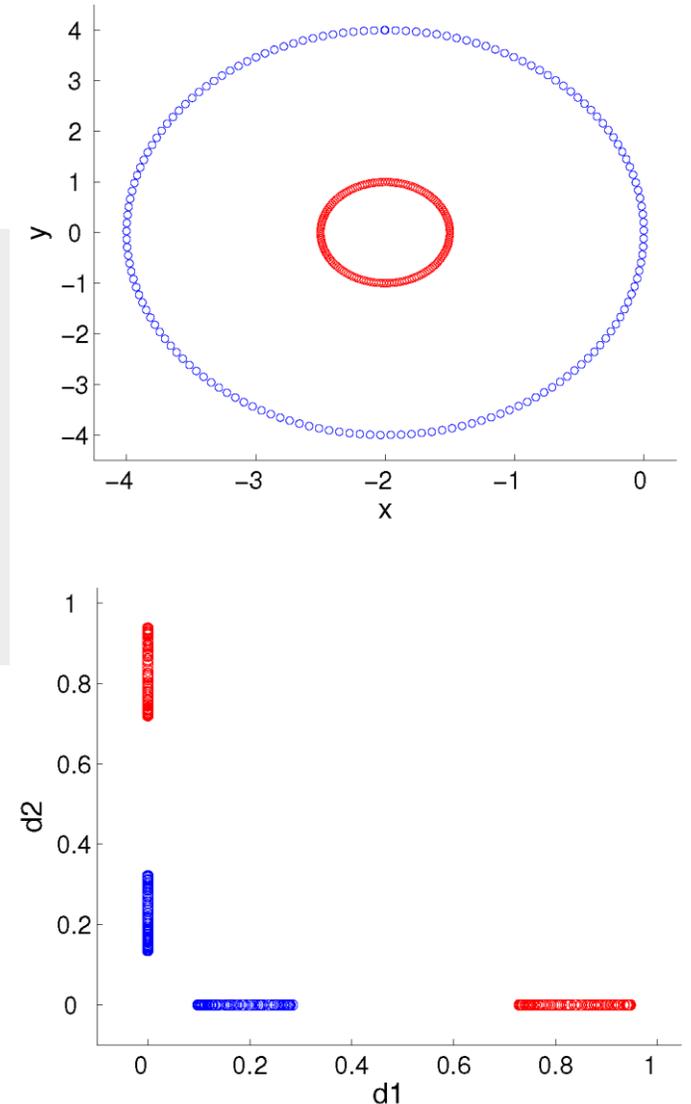
S3 is the most appropriate scheme.

Nonlinear dictionary learning

- Capturing the nonlinear data patterns
- Equiangular kernel learning [Quan et al, CVPR, 2016]

$$\min_{D,C} \|\Phi(Y) - \Phi(D)C\|_F^2, .$$
$$s. t. \|c_i\|_0 \leq k, \forall i, D^T D = I$$

where Φ is a map associated with the kernel K .



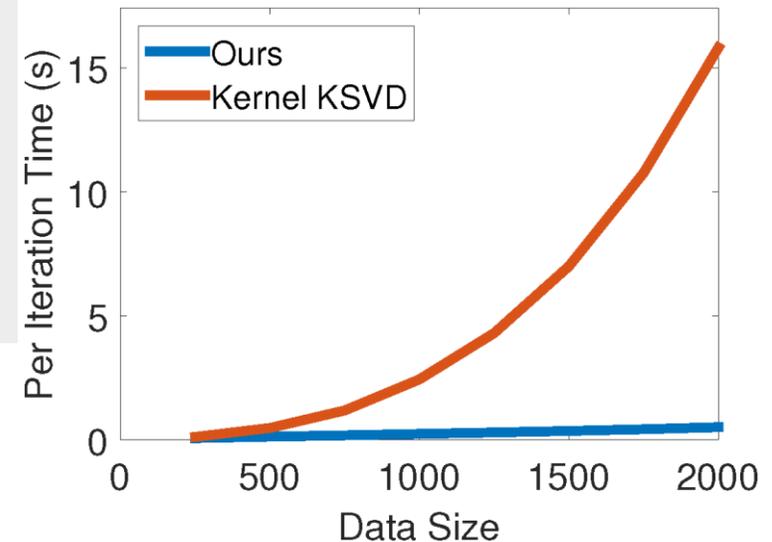
Nonlinear dictionary learning

- Capturing the nonlinear data patterns
- Equiangular kernel learning [Quan et al, CVPR, 2016]

$$\min_{D, C} \|\Phi(Y) - \Phi(D)C\|_F^2, .$$
$$s. t. \|c_i\|_0 \leq k, \forall i, D^T D = I$$

where Φ is a map associated with the kernel K .

- Properties:
 1. $\langle \Phi(d_m), \Phi(d_n) \rangle = \mu_0, \forall m \neq n$ if $K = \psi(\|d_m - d_n\|^2)$.
 2. More scalable than kernel K-SVD (free of $K(Y, Y)$)



Dynamic texture classification

- Recognizing the moving textures with certain stationary temporal changes
- Nonlinear data patterns



Fountain



Flag



Candle

Experiments

- UCLA-DT
 - 50 categories, 200 videos
- DynTex
 - 10 categories, 275 videos
- DynTex++
 - 36 categories, 3600 videos



Sample video

Dataset	DFS	DFS+	LBP-TOP	KGDL	OTF	Ours
UCLA	97.5	97.5	N.A.	N.A.	97.2	98.6
DynTex	74.5	74.8	72.0	75.1	73.5	75.6
DynTex++	89.9	91.7	89.2	92.8	89.8	93.4

Classification accuracy (%)

Discussion

- The hybrid proximal alternating scheme has certain convergence guarantee for multi-block nonconvex problems
- Application based alternating scheme can lead to more efficient and effective algorithm

Future work:

1. Find better initialization
2. Connection between hierarchical dictionary learning and convolutional neural networks

Thank you!