

# Sparse Modeling and Sparse Optimization

Zhouwang Yang

University of Science and Technology of China

April 25, 2021

# Outline I

- ① Compressed Sensing
- ② Sparse Modeling
  - Sparsity-seeking representations
  - Numerical optimization
  - Applications
- ③ Sparse Optimization
  - Sparse Optimization Models
  - Sparse Optimization Algorithms
- ④ Decoupling Noises and Features via  $\ell_1$ -analysis Compressed Sensing
  - Discrete Laplacian regularization smoothing
  - Feature recovering via  $\ell_1$ -analysis optimization

# Outline II

- Experimental results

5 Construction of Manifolds via Consistent Sparse Representations

6 Sparse Representation with Parameterization Optimization

7 Conclusion and Future Work

# Outline I

- ① Compressed Sensing
- ② Sparse Modeling
  - Sparsity-seeking representations
  - Numerical optimization
  - Applications
- ③ Sparse Optimization
  - Sparse Optimization Models
  - Sparse Optimization Algorithms
- ④ Decoupling Noises and Features via  $\ell_1$ -analysis Compressed Sensing
  - Discrete Laplacian regularization smoothing
  - Feature recovering via  $\ell_1$ -analysis optimization

# Outline II

- Experimental results

5 Construction of Manifolds via Consistent Sparse Representations

6 Sparse Representation with Parameterization Optimization

7 Conclusion and Future Work

# Big data era

Electronic commerce data

Social network data

Financial data

Multimedia data

Bioinformatics data

Geometry data

...

# High dimensional data and sparsity

Techniques:

Statistics (Bayesian/Lasso)

Optimization (OMP/BP)

Priors and Transforms

Sparse and Redundant Representations

Low Rank Representations

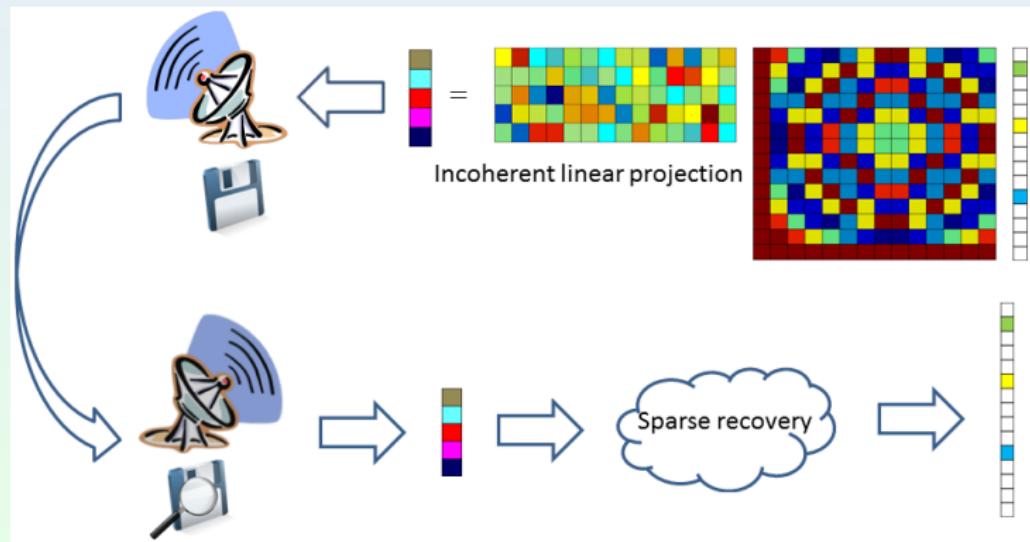
...

# Compressed Sensing

In recent years, Compressed Sensing (CS) has attracted considerable attention in areas of applied mathematics, computer science, and signal processing [Candes and Tao 2005; Donoho 2006; Bruckstein et al. 2009].

## Compressed Sensing

The central insight of CS is that many signals are sparse, i.e., represented using only a few non-zero coefficients in a suitable basis or dictionary and such signals can be recovered from very few measurements (undersampled data) by an optimization algorithm.



# The Sparsest Solution of $Ax = b$

$$(P_0) \quad \min_{\mathbf{x}} \|\mathbf{x}\|_0 \quad \text{s.t.} \quad A\mathbf{x} = \mathbf{b}. \quad (1)$$

For the underdetermined linear system of equations  $Ax = b$  (a full-rank matrix  $A \in \mathbb{R}^{m \times n}$  with  $m \ll n$ ), the following questions are posed:

- Q1: When can uniqueness of the sparsest solution be claimed?
- Q2: Can a candidate solution be tested to verify its (global) optimality?
- Q3: Can the solution be reliably and efficiently found in practice?
- Q4: What performance guarantees can be given for various approximate and practical solvers?

# Uniqueness via the Spark

DEFINITION 01. The spark of a given matrix  $A$  is the smallest number of columns from  $A$  that are linearly dependent.

THEOREM 02. If a system of linear equations  $Ax = b$  has a solution  $x$  obeying  $\|x\|_0 < \text{spark}(A)/2$ , this solution is necessarily the sparsest possible.

## Uniqueness via the Mutual Coherence

DEFINITION 03. The mutual coherence of a given matrix  $A$  is the largest absolute normalized inner product between different columns from  $A$ . Denoting the  $k$ -th column in  $A$  by  $a_k$ , the mutual coherence is given by

$$\mu(A) = \max_{1 \leq i \neq j \leq n} \frac{|a_i^T a_j|}{\|a_i\|_2 \|a_j\|_2}.$$

LEMMA 04. For any matrix  $A \in \mathbb{R}^{m \times n}$ , the following relationship holds:

$$\text{spark}(A) \geq 1 + \frac{1}{\mu(A)}.$$

THEOREM 05. If a system of linear equations  $Ax = b$  has a solution  $x$  obeying  $\|x\|_0 < (1 + 1/\mu(A))/2$ , this solution is necessarily the sparsest possible.

# Pursuit Algorithms

Greedy strategies are usually adopted in solving the problem  $(P_0)$ .

The following algorithm is known in the literature of signal processing by the name *Orthogonal Matching Pursuit* (OMP).

**Task:** Approximate the solution of  $(P_0)$ :  $\min_{\mathbf{x}} \|\mathbf{x}\|_0$  subject to  $\mathbf{A}\mathbf{x} = \mathbf{b}$ .

**Parameters:** We are given the matrix  $\mathbf{A}$ , the vector  $\mathbf{b}$ , and the error threshold  $\epsilon_0$ .

**Initialization:** Initialize  $k = 0$ , and set

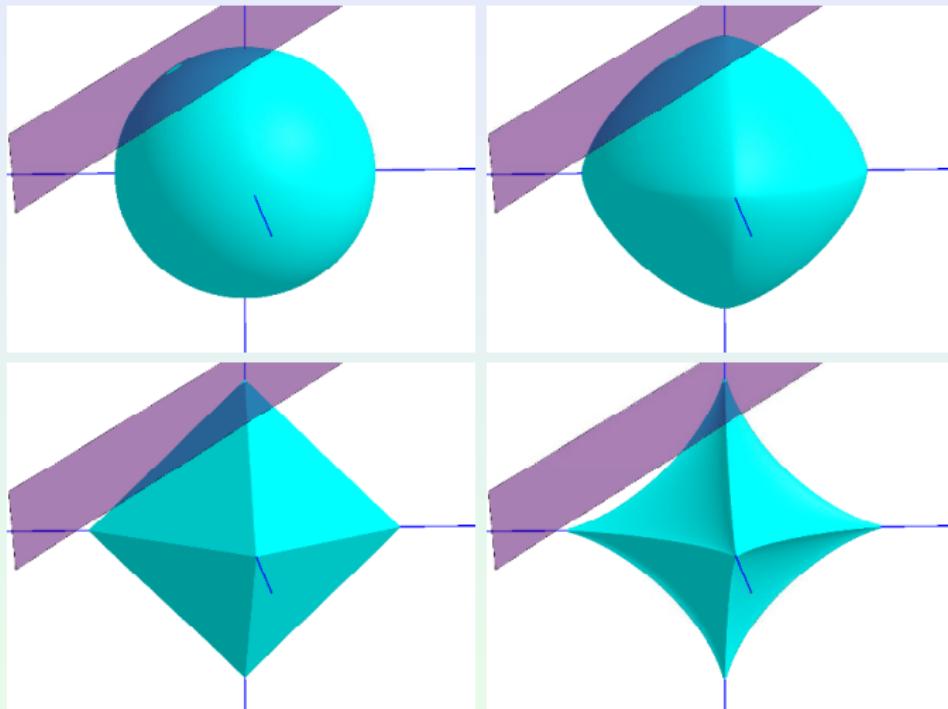
- The initial solution  $\mathbf{x}^0 = 0$ .
- The initial residual  $\mathbf{r}^0 = \mathbf{b} - \mathbf{A}\mathbf{x}^0 = \mathbf{b}$ .
- The initial solution support  $\mathcal{S}^0 = \text{Support}\{\mathbf{x}^0\} = \emptyset$ .

**Main Iteration:** Increment  $k$  by 1 and perform the following steps:

- **Sweep:** Compute the errors  $\epsilon(j) = \min_{z_j} \|\mathbf{a}_j z_j - \mathbf{r}^{k-1}\|_2^2$  for all  $j$  using the optimal choice  $z_j^* = \mathbf{a}_j^T \mathbf{r}^{k-1} / \|\mathbf{a}_j\|_2^2$ .
- **Update Support:** Find a minimizer  $j_0$  of  $\epsilon(j)$ :  $\forall j \notin \mathcal{S}^{k-1}$ ,  $\epsilon(j_0) \leq \epsilon(j)$ , and update  $\mathcal{S}^k = \mathcal{S}^{k-1} \cup \{j_0\}$ .
- **Update Provisional Solution:** Compute  $\mathbf{x}^k$ , the minimizer of  $\|\mathbf{A}\mathbf{x} - \mathbf{b}\|_2^2$  subject to  $\text{Support}\{\mathbf{x}\} = \mathcal{S}^k$ .
- **Update Residual:** Compute  $\mathbf{r}^k = \mathbf{b} - \mathbf{A}\mathbf{x}^k$ .
- **Stopping Rule:** If  $\|\mathbf{r}^k\|_2 < \epsilon_0$ , stop. Otherwise, apply another iteration.

**Output:** The proposed solution is  $\mathbf{x}^k$  obtained after  $k$  iterations.

# Geometry of $\ell_p$ -Norm



# Pursuit Algorithms

Convex relaxation technique is a second way to render ( $P_0$ ) more tractable.

Convexifying with the  $\ell_1$  norm, we come to the new optimization problem

$$(P_1) \quad \min_x \|Wx\|_1 \quad \text{s.t.} \quad Ax = b \quad (2)$$

where  $W$  is a diagonal positive-definite matrix that introduces the precompensating weights.

It was named *Basis Pursuit* (BP) when all the columns of  $A$  are normalized (and thus  $W = I$ ).

# Pursuit Algorithms

**THEOREM 06.** For a system of linear equations  $Ax = b$ , if a solution  $x$  exists obeying  $\|x\|_0 < (1 + 1/\mu(A))/2$ , then an OMP algorithm run with threshold parameter  $\epsilon_0 = 0$  is guaranteed to find it exactly.

**THEOREM 07.** For a system of linear equations  $Ax = b$ , if a solution  $x$  exists obeying  $\|x\|_0 < (1 + 1/\mu(A))/2$ , that solution is both the unique solution of  $(P_1)$  and the unique solution of  $(P_0)$ .

# From Exact to Approximate Solutions

An error-tolerant version of  $(P_0)$  is defined by

$$(P_0^\epsilon) \quad \min_x \|x\|_0 \quad \text{s.t.} \quad \|b - Ax\| \leq \epsilon. \quad (3)$$

**THEOREM 08.** Consider the instance of problem  $(P_0^\epsilon)$  defined by the triplet  $(A; b; \epsilon)$ . Suppose that a sparse vector  $x_0$  satisfies the sparsity constraint  $\|x_0\|_0 < (1 + 1/\mu(A))/2$ , and gives a representation of  $b$  to within error tolerance  $\epsilon$  (i.e.,  $\|b - Ax_0\| \leq \epsilon$ ). Every solution  $x_0^\epsilon$  of  $(P_0^\epsilon)$  must obey

$$\|x_0^\epsilon - x_0\|_2^2 \leq \frac{4\epsilon^2}{1 - \mu(A)(2\|x_0\|_0 - 1)}.$$

# From Exact to Approximate Solutions

An error-tolerant version of  $(P_1)$  is defined by

$$(P_1^\epsilon) \quad \min_{\mathbf{x}} \|\mathbf{x}\|_1 \quad \text{s.t.} \quad \|\mathbf{b} - A\mathbf{x}\| \leq \epsilon. \quad (4)$$

**THEOREM 09.** Consider the instance of problem  $(P_1^\epsilon)$  defined by the triplet  $(A; \mathbf{b}; \epsilon)$ . Suppose that a sparse vector  $\mathbf{x}_0$  is a feasible solution to  $(P_1^\epsilon)$  satisfying the sparsity constraint  $\|\mathbf{x}_0\|_0 < (1 + 1/\mu(A))/4$ . The solution  $\mathbf{x}_1^\epsilon$  of  $(P_1^\epsilon)$  must obey

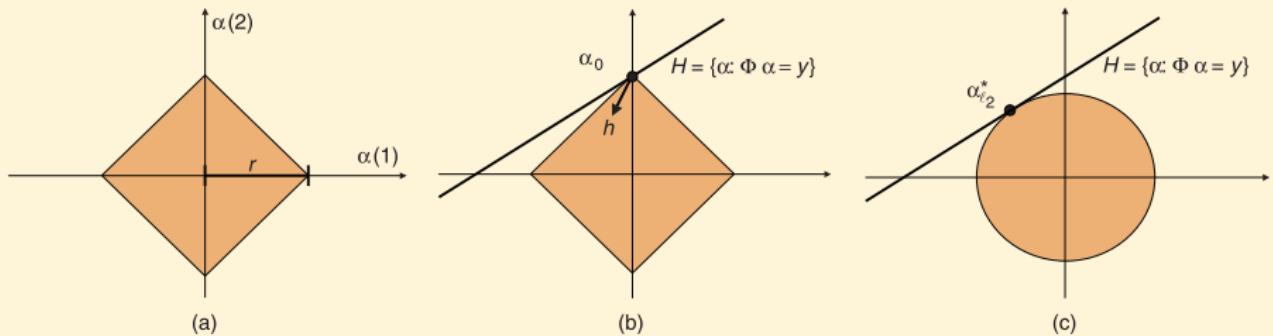
$$\|\mathbf{x}_1^\epsilon - \mathbf{x}_0\|_2^2 \leq \frac{4\epsilon^2}{1 - \mu(A)(4\|\mathbf{x}_0\|_0 - 1)}.$$

# Restricted Isometry Property

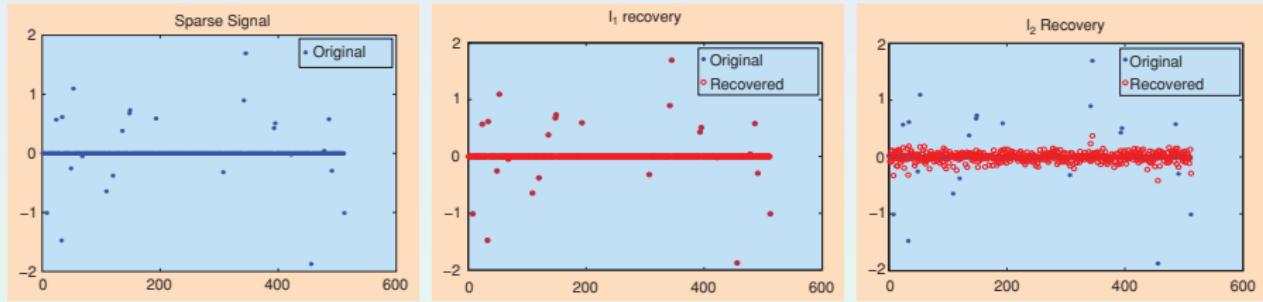
DEFINITION 10. A matrix  $A \in \mathbb{R}^{m \times n}$  is said to have the restricted isometry property  $RIP(\delta; s)$  if each submatrix  $A_s$  formed by combining at most  $s$  columns of  $A$  has its nonzero singular values bounded above by  $1 + \delta$  and below by  $1 - \delta$ .

THEOREM 11. Candès and Tao have shown that  $A \in RIP(\sqrt{2} - 1; 2s)$  implies that  $(P_1)$  and  $(P_0)$  have identical solutions on all  $s$ -sparse vectors and, moreover, that  $(P_1^\epsilon)$  stably approximates the sparsest near-solution of  $b = Ax + v$  with a reasonable stability coefficient.

# $\ell_1$ and $\ell_2$ Recovery



## $\ell_1$ and $\ell_2$ Recovery



# Compressed Sensing

Designing measurement/sensing matrices with favorable properties and constructing suitable transforms/dictionaries are the important research topics in Compressed Sensing.

# Outline I

## ① Compressed Sensing

## ② Sparse Modeling

- Sparsity-seeking representations
- Numerical optimization
- Applications

## ③ Sparse Optimization

- Sparse Optimization Models
- Sparse Optimization Algorithms

## ④ Decoupling Noises and Features via $\ell_1$ -analysis Compressed Sensing

- Discrete Laplacian regularization smoothing
- Feature recovering via  $\ell_1$ -analysis optimization

# Outline II

- Experimental results

5 Construction of Manifolds via Consistent Sparse Representations

6 Sparse Representation with Parameterization Optimization

7 Conclusion and Future Work

# Sparse Modeling

All the previous theorems have shown us that the problem of finding a sparse solution to an under-determined linear system can be given a meaningful definition and can also be computationally tractable.

We now turn to discuss the applicability of these ideas to signal, image, and geometric processing.

# Priors and transforms for signals

The Bayesian framework imposes a Probability-Density-Function (PDF) on the signals – a prior distribution  $P(y)$ .

Priors are extensively used in signal processing, serving in inverse problems, compression, anomaly detection, and more.

## Priors and transforms for signals

Consider the denoising problem: a given image  $b$  is known to be a noisy version of a clean image  $y$ , contaminated by an additive perturbation vector  $v$ , known to have a finite energy  $\|v\|_2 \leq \epsilon$ , i.e.,  $b = y + v$ .

The optimization problem

$$\max_y P(y) \quad \text{s.t.} \quad \|y - b\|_2 \leq \epsilon$$

leads to the most probable image  $\hat{y}$  that is an effective estimate of  $y$ .

This way the prior is exploited for solving the noise cleaning problem. The above formulation of the denoising problem is in fact that Maximum-A-posteriori-Probability (MAP) estimator.

## Priors and transforms for signals

Much effort has been allocated in the signal and image processing communities for forming priors as closed-form expressions.

One very common way to construct  $P(y)$  is to guess its structure based on intuitive expectations from the data content. For example, the Gibbs distribution  $P(y) = \text{Const} \cdot \exp\{-\lambda \|Ly\|_2^2\}$  uses a Laplacian matrix to give an evaluation of the probability of the image  $y$ .

In such a prior, smoothness, measured by the Laplacian operator, is used for judging the probability of the signal.

## Priors and transforms for signals

This prior is well-known and extensively used in signal processing, and is known to be related to both Tikhonov regularization and Wiener filtering.

The prior leads to an optimization problem of the form

$$\min \|Ly\|_2^2 \quad \text{s.t.} \quad \|y - b\|_2 \leq \epsilon$$

which can be converted to

$$\min \lambda \|Ly\|_2^2 + \|y - b\|_2^2$$

where we have replaced the constraint by an equivalent penalty.

The closed-form solution is easily obtained as

$$\hat{y} = (I + \lambda L^T L)^{-1} b.$$

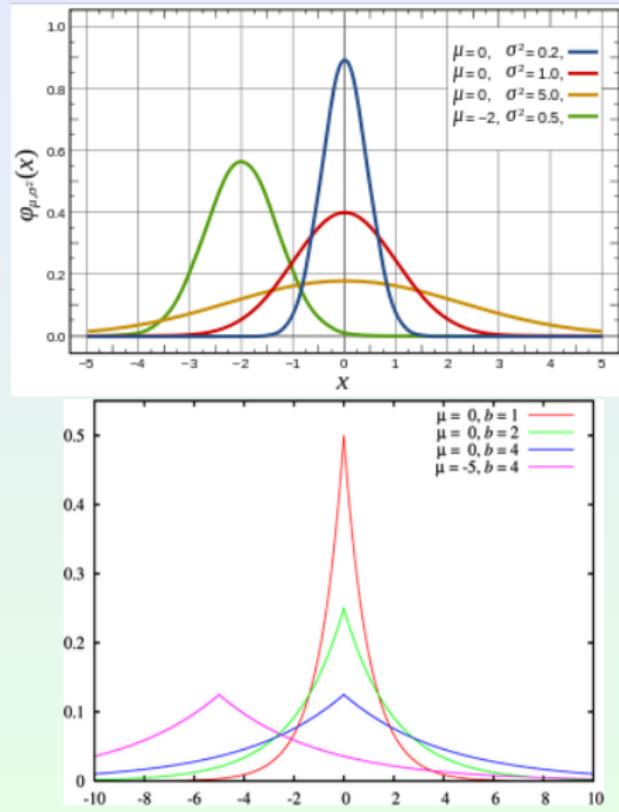
## Priors and transforms for signals

The above specific prior stressing smoothness is known to cause blurring of the image when used in various restoration tasks. The remedy for this problem was found to be the replacement of the  $\ell_2$ -norm by a more robust measure, such as an  $\ell_1$ -norm, that allows heavy tails for the distribution of the values of  $Ly$ .

A prior of the form  $P(y) = Const \cdot \exp\{-\lambda \|Ly\|_1\}$  is far more versatile and thus became popular in recent years.

Similar to this option is the Total-Variation (TV) prior  $P(y) = Const \cdot \exp\{-\lambda \|y\|_{TV}\}$  [Rudin, Osher, and Fatemi, 1993] that also promotes smoothness, but differently, by replacing the Laplacian with gradient norms.

# Priors and transforms for signals



## Priors and transforms for signals

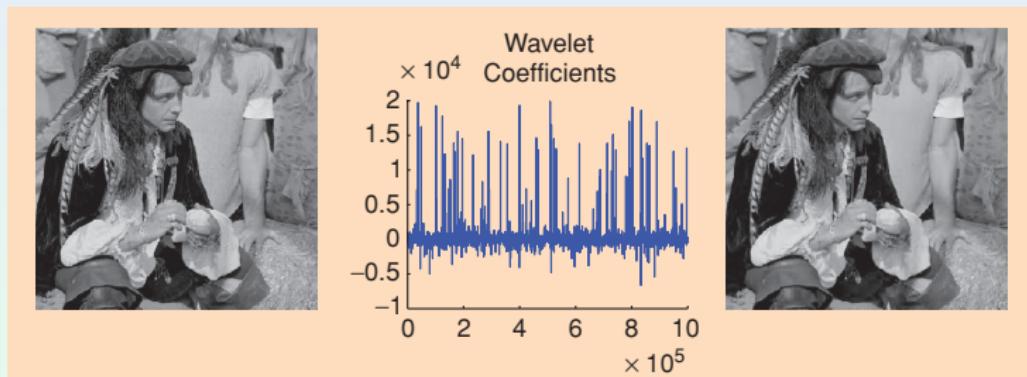
A different property that can be used for constructing a prior is assuming a structure on the signals transform-coefficients.

One such example is the JPEG compression algorithm, which relies on the fact that 2D-DCT coefficients of small image patches tend to behave in a predicted way (being concentrated around the origin).

Another well-known example refers to the wavelet transform of signals and images, where the coefficients are expected to be sparse, most of them tending to zero while few remain active.

# Priors and transforms for signals

For a signal  $y$ , the wavelet transform is given by  $Ty$  where the matrix  $T$  is a specially designed orthogonal matrix that contains in its rows spatial derivatives of varying scale, thereby providing what is known as multi-scale analysis of the signal.



Therefore, the prior in this case becomes

$$P(y) = \text{Const} \cdot \exp\{-\lambda \|Ty\|_p^p\} \text{ with } p \leq 1 \text{ to promote sparsity.}$$

# Priors and transforms for signals

A rich family of signal priors assign likelihood for an image based on the behavior of its transform coefficients  $Ty$ . In the signal and image processing literature, such priors were postulated in conjunction with a variety of transforms, such as

- the Discrete-Fourier-Transform (DFT)
- the Discrete-Cosine-Transform (DCT)
- the Hadamard-Transform (HT)
- the Principal-Component-Analysis (PCA)

# Bayesian perspective

- Using Bayes' rule, the posterior probability of  $y$  given the measurements is formulated by

$$P(y | z) = \frac{P(z | y)P(y)}{P(z)}.$$

- Considering the fact that the denominator  $P(z)$  is not a function of the unknown  $y$ , and as such it can be disregarded, the MAP estimation amounts to

$$\hat{y}_{\text{MAP}} = \arg \max_y P(y | z) = \arg \max_y P(z | y)P(y).$$

- The probability  $P(z | y)$  is known as the likelihood function, and the probability  $P(y)$  is the known/unknown's prior.

# The Sparse-Land model

The linear system  $Dx = y$  can be interpreted as a way of constructing signals  $y$ . Every column in  $D$  is a possible signal in  $\mathbb{R}^n$  – we refer to these  $m$  columns as atomic signals, and the matrix  $D$  displays a dictionary of atoms.

The multiplication of  $D$  by a sparse vector  $x$  with  $\|x\|_0 = k_0 \ll n$  produces a linear combination of  $k_0$  atoms with varying portions, generating the signal  $y$ . The vector  $x$  that generates  $y$  will be called its representation.

# The Sparse-Land model

The Sparse-Land model  $\mathcal{M}(D, k_0, \alpha, \epsilon)$ :

$$y = Dx + v$$

- Consider all the possible sparse representation vectors with cardinality  $\|x\|_0 = k_0 \ll n$ , and assume that this set of  $C_m^{k_0}$  possible cardinalities are drawn with uniform probability.
- Assume further that the non-zero entries in  $x$  are drawn from the zero-mean Gaussian distribution  $Const \cdot \exp\{-\alpha x_i^2\}$ .
- Postulate that the observations are contaminated by a random perturbation (noise) vector  $v \in \mathbb{R}^n$  with bounded power  $\|v\|_2 \leq \epsilon$ .

# The quest for a dictionary

In the quest for the proper dictionary to use in applications, one line of work considers choosing pre-constructed dictionaries, such as undecimated wavelets, steerable wavelets, contourlets, curvelets, and more.

Some of these proposed dictionaries (which are often referred to also as transforms) are accompanied by a detailed theoretical analysis establishing the sparsity of the representation coefficients for such simplified content of signals.

## The quest for a dictionary

While pre-constructed or adapted dictionaries typically lead to fast transforms, they are typically limited in their ability to sparsify the signals they are designed to handle. Furthermore, most of those dictionaries are restricted to signals/images of a certain type, and cannot be used for a new and arbitrary family of signals of interest.

This leads us to yet another approach for obtaining dictionaries that overcomes these limitations – by adopting a learning point-of-view.

As opposed to the pre-constructed and adapted dictionaries, the learning method is able to adapt to any family of signals that complies with the Sparse-Land model.

# Dictionary learning

Assume that a training database  $\{y_j\}_{j=1}^N$  is given, and thought to have been generated by some fixed but unknown model  $\mathcal{M}(D, k_0, \alpha, \epsilon)$ .

- Control the deviation:

$$\min_{D, \{x_j\}_{j=1}^N} \sum_{j=1}^N \|x_j\|_0 \quad \text{s.t.} \quad \|y_j - Dx_j\|_2 \leq \epsilon, \quad j = 1, \dots, N$$

- Control the sparsity:

$$\min_{D, \{x_j\}_{j=1}^N} \sum_{j=1}^N \|y_j - Dx_j\|_2^2 \quad \text{s.t.} \quad \|x_j\|_0 \leq k_0, \quad j = 1, \dots, N$$

# Analysis versus Synthesis

- Synthesis based modeling

$$(P_s) \quad \hat{y}_s = D \cdot \arg \min_x \|x\|_p \quad \text{s.t.} \quad \|z - Dx\| \leq \epsilon. \quad (5)$$

- Analysis based modeling

$$(P_a) \quad \hat{y}_a = \arg \min_y \|Ty\|_p \quad \text{s.t.} \quad \|z - y\| \leq \epsilon. \quad (6)$$

# Dictionary learning algorithms

The optimization problem for dictionary learning (sparse representations and coding):

$$\min_{D, X} \|Y - DX\|_{Frob} \quad \text{s.t.} \quad \|x_j\|_0 \leq k_0, \quad j = 1, \dots, N \quad (7)$$

where

$$Y = (y_1, \dots, y_N) \in \mathbb{R}^{n \times N},$$

$$D = (d_1, \dots, d_m) \in \mathbb{R}^{n \times m},$$

$$X = (x_1, \dots, x_N) \in \mathbb{R}^{m \times N}.$$

# Dictionary learning algorithms

There are two training mechanisms, the first named Method of Optimal Directions (MOD) by Engan et al., and the second named K-SVD, by Aharon et al..

- MOD
- K-SVD
- .....

# Applications

- Image deblurring
- Facial image compression
- Image denoising
- Image inpainting
- .....

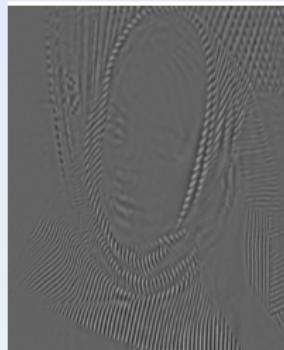
# Applications



Original image



Cartoon part



Texture part

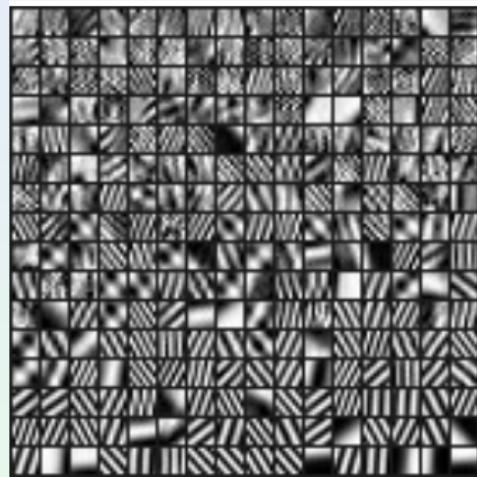
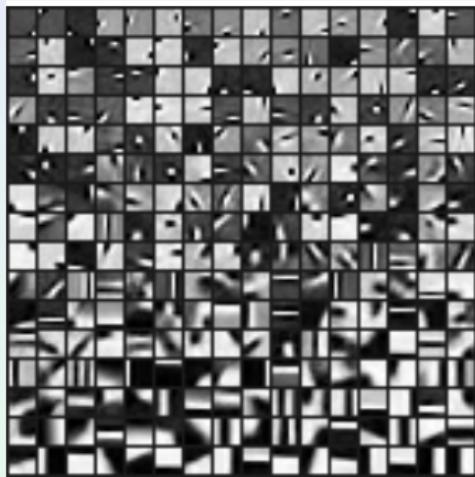
Image *inpainting* [2, 10, 20, 38] is the process of inserting data in a designated region of a still or moving image. Applications range from removing objects from images to restoring damaged paintings and photographs. Inpainting produces a revised image in which the missing data is seamlessly merged into the image in a way that is not detectable by a typical viewer. Traditionally, inpainting has been done by professional artists. For photographs, inpainting is used to reverse deterioration, such as scratches and dust spots in film, or to remove elements (e.g., removal of stamped date from photographs, the infamous “airbrushed” enemies [20]). A current active area of research is



Image with missing data

Inpainting result

# Applications



# Applications

Original Image



Noisy Image (22.1307 dB,  $\sigma=20$ )



Denoised Image Using  
Global Trained Dictionary (28.8528 dB)



Denoised Image Using  
Adaptive Dictionary (30.8295 dB)



## References

1. Emmanuel J. Candès, Terence Tao. Decoding by linear programming. *Information Theory, IEEE Transactions on* 51.12 (2005): 4203-4215.
2. Emmanuel J. Candès, Justin Romberg, Terence Tao. Robust uncertainty principles: Exact signal reconstruction from highly incomplete frequency information. *Information Theory, IEEE Transactions on* 52.2 (2006): 489-509.
3. Emmanuel J. Candès, Michael B. Wakin. An introduction to compressive sampling. *Signal Processing Magazine, IEEE* 25.2 (2008): 21-30.
4. Alfred M. Bruckstein, David L. Donoho, Michael Elad. From sparse solutions of systems of equations to sparse modeling of signals and images. *SIAM review* 51.1 (2009): 34-81.
5. Michael Elad. Sparse and Redundant Representations: From Theory to Applications in Signal and Image Processing. Springer, 2010.

# Outline I

## ① Compressed Sensing

## ② Sparse Modeling

- Sparsity-seeking representations
- Numerical optimization
- Applications

## ③ Sparse Optimization

- Sparse Optimization Models
- Sparse Optimization Algorithms

## ④ Decoupling Noises and Features via $\ell_1$ -analysis Compressed Sensing

- Discrete Laplacian regularization smoothing
- Feature recovering via  $\ell_1$ -analysis optimization

# Outline II

- Experimental results

5 Construction of Manifolds via Consistent Sparse Representations

6 Sparse Representation with Parameterization Optimization

7 Conclusion and Future Work

# Sparse Optimization

Many problems of recent interest in statistics and related areas can be posed in the framework of sparse optimization. Due to the explosion in size and complexity of modern data analysis (BigData), it is increasingly important to be able to solve problems with a very large number of features, training examples, or both.

## 0-norm optimization

$$(P_0) \quad \min_{\mathbf{x}} \|\mathbf{x}\|_0 \quad \text{s.t.} \quad A\mathbf{x} = \mathbf{b}. \quad (8)$$

$$(P_0^\epsilon) \quad \min_{\mathbf{x}} \|\mathbf{x}\|_0 \quad \text{s.t.} \quad \|\mathbf{b} - A\mathbf{x}\| \leq \epsilon. \quad (9)$$

# Greedy algorithms

Greedy strategies are usually adopted in solving the 0-norm problems. The following algorithm is known in the literature of signal processing by the name *Orthogonal Matching Pursuit* (OMP).

**Task:** Approximate the solution of  $(P_0)$ :  $\min_{\mathbf{x}} \|\mathbf{x}\|_0$  subject to  $\mathbf{Ax} = \mathbf{b}$ .

**Parameters:** We are given the matrix  $\mathbf{A}$ , the vector  $\mathbf{b}$ , and the error threshold  $\epsilon_0$ .

**Initialization:** Initialize  $k = 0$ , and set

- The initial solution  $\mathbf{x}^0 = \mathbf{0}$ .
- The initial residual  $\mathbf{r}^0 = \mathbf{b} - \mathbf{Ax}^0 = \mathbf{b}$ .
- The initial solution support  $\mathcal{S}^0 = \text{Support}\{\mathbf{x}^0\} = \emptyset$ .

**Main Iteration:** Increment  $k$  by 1 and perform the following steps:

- **Sweep:** Compute the errors  $\epsilon(j) = \min_{z_j} \|\mathbf{a}_j z_j - \mathbf{r}^{k-1}\|_2^2$  for all  $j$  using the optimal choice  $z_j^* = \mathbf{a}_j^T \mathbf{r}^{k-1} / \|\mathbf{a}_j\|_2^2$ .
- **Update Support:** Find a minimizer  $j_0$  of  $\epsilon(j)$ :  $\forall j \notin \mathcal{S}^{k-1}$ ,  $\epsilon(j_0) \leq \epsilon(j)$ , and update  $\mathcal{S}^k = \mathcal{S}^{k-1} \cup \{j_0\}$ .
- **Update Provisional Solution:** Compute  $\mathbf{x}^k$ , the minimizer of  $\|\mathbf{Ax} - \mathbf{b}\|_2^2$  subject to  $\text{Support}\{\mathbf{x}\} = \mathcal{S}^k$ .
- **Update Residual:** Compute  $\mathbf{r}^k = \mathbf{b} - \mathbf{Ax}^k$ .
- **Stopping Rule:** If  $\|\mathbf{r}^k\|_2 < \epsilon_0$ , stop. Otherwise, apply another iteration.

**Output:** The proposed solution is  $\mathbf{x}^k$  obtained after  $k$  iterations.

# Dictionary learning

The optimization model of dictionary learning for sparse and redundant representations:

$$\min_{D, X} \|Y - DX\|_{Frob} \quad \text{s.t.} \quad \|x_j\|_0 \leq k_0, \quad j = 1, \dots, N \quad (10)$$

where

$$Y = (y_1, \dots, y_N) \in \mathbb{R}^{n \times N},$$

$$D = (d_1, \dots, d_m) \in \mathbb{R}^{n \times m},$$

$$X = (x_1, \dots, x_N) \in \mathbb{R}^{m \times N}.$$

# Dictionary learning

There are two training mechanisms, the first named Method of Optimal Directions (MOD) by Engan et al., and the second named K-SVD, by Aharon et al..

- MOD
- K-SVD
- .....

## Convex relaxation

Convex relaxation technique is a way to render 0-norm more tractable.

Convexifying with the  $\ell_1$  norm, we come to the new optimization problem

$$(P_1) \quad \min_x \|Wx\|_1 \quad \text{s.t.} \quad Ax = b \quad (11)$$

where  $W$  is a diagonal positive-definite matrix that introduces the precompensating weights.

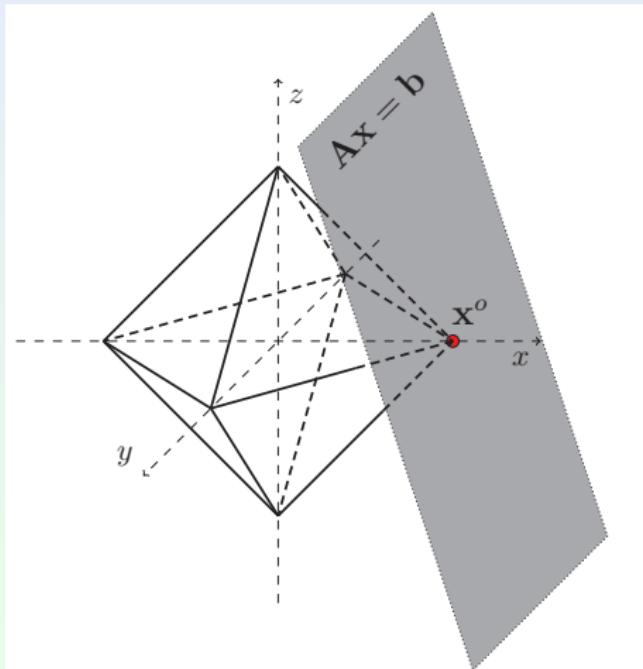
An error-tolerant version of  $(P_1)$  is defined by

$$(P_1^\epsilon) \quad \min_x \|Wx\|_1 \quad \text{s.t.} \quad \|b - Ax\| \leq \epsilon. \quad (12)$$

It was named *Basis Pursuit* (BP) when all the columns of  $A$  are normalized (and thus  $W = I$ ).

## Basic pursuit

$$(BP) \quad \min_x \|x\|_1 \quad \text{s.t.} \quad Ax = b.$$



# BP denoising and LASSO

$$(BP_\tau) \quad \min_{\mathbf{x}} \|\mathbf{A}\mathbf{x} - \mathbf{b}\|_2^2 \quad \text{s.t.} \quad \|\mathbf{x}\|_1 \leq \tau,$$

$$(BP_\mu) \quad \min_{\mathbf{x}} \|\mathbf{x}\|_1 + \frac{\mu}{2} \|\mathbf{A}\mathbf{x} - \mathbf{b}\|_2^2,$$

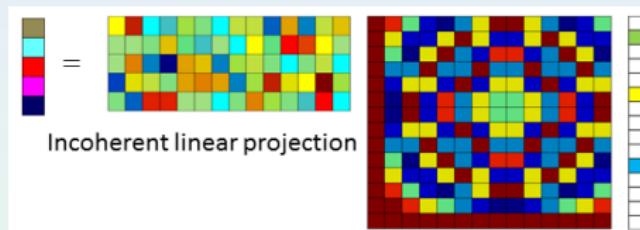
$$(BP_\delta) \quad \min_{\mathbf{x}} \|\mathbf{x}\|_1 \quad \text{s.t.} \quad \|\mathbf{A}\mathbf{x} - \mathbf{b}\|_2 \leq \delta.$$

Questions:

- Are they equivalent? and in what sense?
- How to choose parameters?

# Sparse under basis $\Psi$

$$\min_s \{ \|s\|_1 : A\Psi s = b \}$$



If  $\Psi$  is orthogonal, the problem is equivalent to

$$\min_x \{ \|\Psi^*x\|_1 : Ax = b \}.$$

# Sparse after transform $\mathcal{L}$

$$\min_{\mathbf{x}} \{\|\mathcal{L}\mathbf{x}\|_1 : A\mathbf{x} = \mathbf{b}\}$$

Examples of  $\mathcal{L}$ :

- DCT, wavelets, curvelets, ridgelets, ...
- tight frames, Gabor, ...
- total (generalized) variation

Ref: E. J. Cands, Y. Eldar, D. Needell and P. Randall. Compressed sensing with coherent and redundant dictionaries. Applied and Computational Harmonic Analysis, 31(1): 59-73.

## Joint/group sparsity

Decompose  $\{1, 2, \dots, n\} = \mathcal{G}_1 \cup \mathcal{G}_2 \cup \dots \cup \mathcal{G}_S$ , and  $\mathcal{G}_i \cap \mathcal{G}_j = \emptyset, i \neq j$ .

Joint/group sparse recovery model:

$$\min_{\mathbf{x}} \{\|\mathbf{x}\|_{\mathcal{G},2,1} : A\mathbf{x} = \mathbf{b}\}$$

where

$$\|\mathbf{x}\|_{\mathcal{G},2,1} = \sum_{s=1}^S w_s \|\mathbf{x}_{\mathcal{G}_s}\|_2.$$

## Side constraints

- Nonnegativity:  $x \geq 0$
- Box constraints:  $lb \leq x \leq ub$
- Linear inequalities:  $Qx \leq c$

They generate “corners” and can be very effective in practice.

# Shrinkage

- Shrinkage is popular in sparse optimization algorithms
- In optimization, non-smooth functions like  $\ell_1$  has difficulty using general smooth optimization methods.
- But,  $\ell_1$  is component-wise separable, so it does get along well with separable (smooth or non-smooth) functions.
- For example,

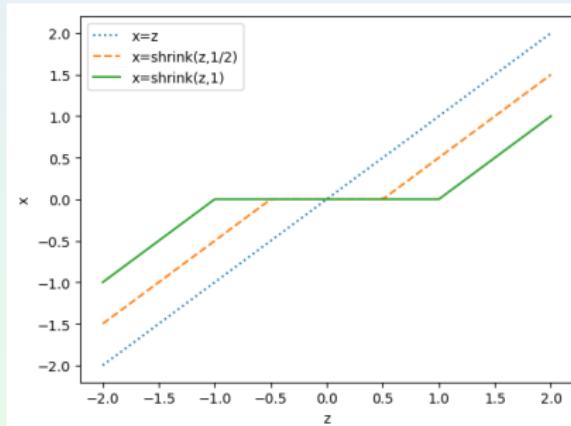
$$\min_{\mathbf{x}} \|\mathbf{x}\|_1 + \frac{1}{2\tau} \|\mathbf{x} - \mathbf{z}\|_2^2$$

is equivalent to solving  $\min_{x_i} |x_i| + \frac{1}{2\tau} |x_i - z_i|^2$  over each  $i$ .

# Soft-thresholding shrinkage

The problem is separable and has an explicit solution

$$(\text{shrink}(z, \tau))_i = \begin{cases} z_i - \tau & z_i > \tau, \\ 0 & -\tau \leq z_i \leq \tau, \\ z_i + \tau & z_i < -\tau. \end{cases}$$



The shrinkage operator can be written in Matlab code as:

$$\mathbf{x} = \max(\text{abs}(\mathbf{z}) - \tau, 0) * \text{sign}(\mathbf{z}).$$

## Soft-thresholding shrinkage

- The following problem is called Moreau-Yosida regularization

$$\min_x r(x) + \frac{1}{2\tau} \|x - z\|_2^2.$$

- For example  $r(x) = \|x\|_2$ , the solution to

$$\min_x \|x\|_2 + \frac{1}{2\tau} \|x - z\|_2^2$$

is, if we treat  $0/0 = 0$ ,

$$x_{opt} = \max\{\|z\|_2 - \tau, 0\} \cdot (z/\|z\|_2).$$

- Used in joint/group-sparse recovery algorithms.

## Soft-thresholding shrinkage

- Consider the following nuclear norm optimization

$$\min_{\mathbf{X}} \|\mathbf{X}\|_* + \frac{1}{2\tau} \|\mathbf{X} - \mathbf{Z}\|_F^2.$$

Let  $\mathbf{Z} = \mathbf{U}\Sigma\mathbf{V}^T$  be the singular value decomposition of  $\mathbf{Z}$ .

- Let  $\hat{\Sigma}$  be the diagonal matrix with diagonal entries

$$\text{diag}(\hat{\Sigma}) = \text{shrink}(\text{diag}(\Sigma), \tau),$$

then

$$\mathbf{X}_{opt} = \mathbf{U}\hat{\Sigma}\mathbf{V}^T.$$

- In general, matrix problems with only unitary-invariant functions (e.g.,  $\|\cdot\|_*$ ,  $\|\cdot\|_F$ , spectral norm, trace) and constraints (e.g., positive or negative semi-definiteness) typically reduce to vector problems regarding singular values.

# Prox-linear algorithm

Consider the general form

$$\min_x r(x) + f(x).$$

where  $r$  is the regularization function and  $f$  is the data fidelity function.

The prox-linear algorithm is:

$$x^{k+1} = \arg \min_x r(x) + f(x^k) + \langle \nabla f(x^k), x - x^k \rangle + \frac{1}{2\delta_k} \|x - x^k\|_2^2.$$

The last term keeps  $x^{k+1}$  close to  $x^k$ , and the parameter  $\delta_k$  determines the step size. It is equivalent to

$$x^{k+1} = \arg \min_x r(x) + \frac{1}{2\delta_k} \|x - (x^k - \delta_k \nabla f(x^k))\|_2^2.$$

# Alternating direction method of multipliers (ADMM)

The Alternating Direction Method of Multipliers (ADMM) was developed in the 1970s, with roots in the 1950s, and is equivalent or closely related to many other algorithms, such as dual decomposition, the method of multipliers, Douglas-Rachford splitting, Spingarns method of partial inverses, Dykstras alternating projections, Bregman iterative algorithms for 1-norm problems, proximal methods, and others.

# ADMM

The ADMM can be applied to a wide variety of statistical and machine learning problems of recent interest, including the lasso, sparse logistic regression, basis pursuit, covariance selection, support vector machines, and many others.

# ADMM

$$\min_{X \in C^{n \times T}} \mu \|X\|_p + \|AX - B\|_q \quad (13)$$

Let  $p := \{2, 1\}$ ,  $q := \{1, 1\}$  which denote joint convex norm, we have

$$\min_{X \in C^{n \times T}} \mu \|X\|_{2,1} + \|AX - B\|_{1,1}$$

where  $\|X\|_{2,1} = \sum_{i=1}^n \sqrt{\sum_{j=1}^T x_{ij}^2}$ ,  $\|X\|_{1,1} = \sum_{i=1}^n \sum_{j=1}^T |x_{ij}|$ .

For example  $T = 1$ ,

$$\min_{x \in C^n} \mu \|x\|_p + \|Ax - b\|_q.$$

# ADMM

$$\min_{X \in C^{n \times T}} \mu \|X\|_p + \|AX - B\|_q \quad (13)$$

Let  $p := \{2, 1\}$ ,  $q := \{1, 1\}$  which denote joint convex norm, we have

$$\min_{X \in C^{n \times T}} \mu \|X\|_{2,1} + \|AX - B\|_{1,1}$$

where  $\|X\|_{2,1} = \sum_{i=1}^n \sqrt{\sum_{j=1}^T x_{ij}^2}$ ,  $\|X\|_{1,1} = \sum_{i=1}^n \sum_{j=1}^T |x_{ij}|$ .

For example  $T = 1$ ,

$$\min_{x \in C^n} \mu \|x\|_p + \|Ax - b\|_q.$$

# ADMM

$$\min_{X \in C^{n \times T}} \mu \|X\|_p + \|AX - B\|_q \quad (13)$$

Let  $p := \{2, 1\}$ ,  $q := \{1, 1\}$  which denote joint convex norm, we have

$$\min_{X \in C^{n \times T}} \mu \|X\|_{2,1} + \|AX - B\|_{1,1}$$

where  $\|X\|_{2,1} = \sum_{i=1}^n \sqrt{\sum_{j=1}^T x_{ij}^2}$ ,  $\|X\|_{1,1} = \sum_{i=1}^n \sum_{j=1}^T |x_{ij}|$ .

For example  $T = 1$ ,

$$\min_{x \in C^n} \mu \|x\|_p + \|Ax - b\|_q.$$

# ADMM

$$\begin{aligned} & \min \mu\|z\|_p + \|y\|_q \\ \text{s.t. } & x - z = 0 \\ & Ax - y = b \end{aligned} \tag{14}$$

$$\begin{aligned} L(x, y, z, \lambda_y, \lambda_z, \rho) = & \mu\|z\|_p + \|y\|_q + \operatorname{Re}(\lambda_z^T(x - z) + \lambda_y^T(Ax - y - b)) \\ & + \frac{\rho}{2}(\|x - z\|_2^2 + \|Ax - y - b\|_2^2) \end{aligned} \tag{15}$$

where  $\lambda_y \in C^n$ ,  $\lambda_z \in C^m$  are the Lagrangian multipliers and  $\rho > 0$  is a penalty parameter.

# ADMM

$$\begin{aligned} & \min \mu\|z\|_p + \|y\|_q \\ \text{s.t. } & x - z = 0 \\ & Ax - y = b \end{aligned} \tag{14}$$

$$\begin{aligned} L(x, y, z, \lambda_y, \lambda_z, \rho) = & \mu\|z\|_p + \|y\|_q + \operatorname{Re}(\lambda_z^T(x - z) + \lambda_y^T(Ax - y - b)) \\ & + \frac{\rho}{2}(\|x - z\|_2^2 + \|Ax - y - b\|_2^2) \end{aligned} \tag{15}$$

where  $\lambda_y \in C^n$ ,  $\lambda_z \in C^m$  are the Lagrangian multipliers and  $\rho > 0$  is a penalty parameter.

$$\begin{cases} x^{k+1} := \arg \min \frac{1}{2}(\|x - z^k + u_z^k\|_2^2 + \|Ax - y^k - b + u_y^k\|_2^2) \\ y^{k+1} := \arg \min \|y\|_q + \frac{\rho}{2}\|y - (Ax^{k+1} - b) - u_y^k\|_2^2 \\ z^{k+1} := \arg \min \mu\|z\|_p + \frac{\rho}{2}\|z - x^{k+1} - u_z^k\|_2^2 \end{cases} \quad (16)$$

After solving three subproblems, we update the Lagrangian multipliers as follows:

$$\begin{cases} u_z^{k+1} = u_z^k + \gamma(x^{k+1} - z^{k+1}) \\ u_y^{k+1} = u_y^k + \gamma(Ax^{k+1} - y^{k+1} - b) \end{cases} \quad (17)$$

where  $u_y = \frac{1}{\rho}\lambda_y$ ,  $u_z = \frac{1}{\rho}\lambda_z$ ,  $\gamma > 0$  is the step size.

# Outline I

- ① Compressed Sensing
- ② Sparse Modeling
  - Sparsity-seeking representations
  - Numerical optimization
  - Applications
- ③ Sparse Optimization
  - Sparse Optimization Models
  - Sparse Optimization Algorithms
- ④ Decoupling Noises and Features via  $\ell_1$ -analysis Compressed Sensing
  - Discrete Laplacian regularization smoothing
  - Feature recovering via  $\ell_1$ -analysis optimization

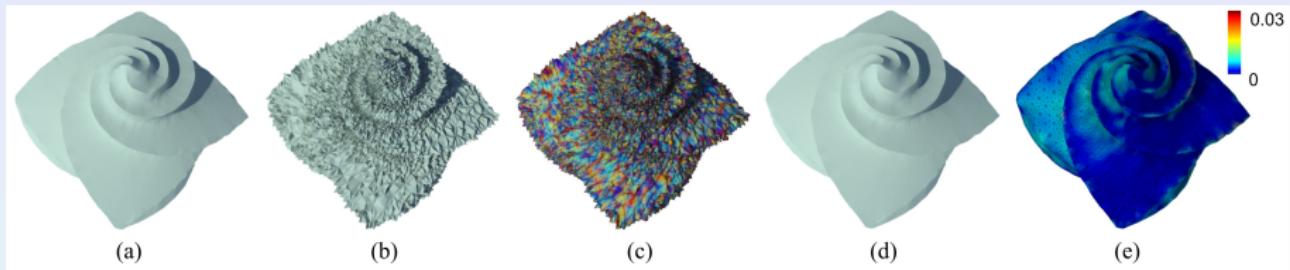
# Outline II

- Experimental results

5 Construction of Manifolds via Consistent Sparse Representations

6 Sparse Representation with Parameterization Optimization

7 Conclusion and Future Work



**Figure:** Our approach is able to faithfully recover sharp features corrupted by noise. (a) The octa-flower model as a ground truth; (b) the model artificially corrupted by independent and identically distributed (i.i.d.) noise (with zero mean and standard deviation  $\sigma = 0.01$  of the diagonal of the bounding box of the model); (c) the error map of (b); (d) the denoising result by our approach; (e) the error map of (d). The colored models (c,e) visualize the errors between the processed models and the ground truth model (a). Note that sharp features such as creases and corners are well retained in our result (d).

# Feature preserving surface denoising

- Surface fairing [Taubin 1995; Desbrun 1999]
- Anisotropic geometric diffusion [Clarenz 2000; Bajaj and Xu 2003]
- Bilateral filtering [Fleishman et al. 2003]
- Non-iterative smoothing [Jones et al. 2003]
- Anisotropic filtering [Hildebrandt and Polthier 2004]
- Normal-based bilateral filtering [Zheng et al. 2010]
- ...

# Statistics

- Splines [Duchon 1977; Utreras 1988; Wahba 1990]
- Nonparametric regression [Stone 1982]
- Generalized cross validation (GCV) [Wahba and Wendelberger 1980]
- ...

## Problem formulation

Assume that  $P = \{p_i\}_1^n$  are sampled possibly with noise from a  $C^2$ -smooth surface  $S$ , i.e.,

$$p_i = s_i + \varepsilon_i n_i, \quad i = 1, \dots, n, \quad (18)$$

where  $n_i$  is the unit normal of surface at  $s_i \in S$  and  $\varepsilon_i, i = 1, \dots, n$ , are i.i.d. random variables with zero mean and finite variance  $\sigma^2$ .

The goal of denoising is to produce a smooth surface  $\hat{S} = \{\hat{s}_i\}$  to approximate the true underlying surface  $S$  as much as possible.

## Problem formulation

To find a  $C^2$ -smooth surface  $\hat{S}$  which is a reasonable estimate of the true underlying surface  $S$ , we formulate it as the following variational minimization problem

$$\hat{S} = \arg \min_S \frac{1}{n} \sum_{i=1}^n d^2(p_i, S) + \lambda \int_S (2H)^2, \quad (19)$$

where  $d(p_i, S)$  is the geometric distance from  $p_i$  to  $S$ ,  $H$  is the mean curvature of  $S$ , and  $\lambda$  is a smoothness parameter.

# Problem formulation

To simplify the computation, we replace the smoothness term by its discrete approximation

$$\begin{aligned}\mathcal{J}(S) &= \frac{1}{n} \sum_{i=1}^n 4H^2(s_i) \\ &= \frac{1}{n} \sum_{i=1}^n \|\Delta_S s_i\|^2 \simeq \frac{1}{n} \sum_{i=1}^n \|L_i S\|^2,\end{aligned}\tag{20}$$

where  $S = (s_1, \dots, s_n)^T \in \mathbb{R}^{n \times 3}$ ,  $\Delta_S$  is the Laplace-Beltrami operator on surface  $S$ , and  $L = (L_1^T, \dots, L_n^T)^T$  is the discrete Laplacian matrix. The second equality in (20) is derived by  $\Delta_S s_i = 2H n_i$ .

## DLRS model

Thus we arrive at the following denoising model:

$$\hat{S} = \arg \min_S \sum_{i=1}^n \|p_i - s_i\|^2 + \lambda \sum_{i=1}^n \|L_i S\|^2. \quad (21)$$

We call it the discrete Laplacian regularization smoothing (DLRS) model.

## DLRS model

Denote  $P = (p_1, \dots, p_n)^T \in \mathbb{R}^{n \times 3}$ . The DLRS model eventually leads to linear systems

$$(I_n + \lambda M) \hat{S} = P,$$

where  $I_n$  is the  $n \times n$  identity matrix and  $M = L^T L = \sum_{i=1}^n L_i^T L_i$ .

Thus, given a specific value of the smoothness parameter  $\lambda$ , we have the solution

$$\hat{S}_n(\lambda) = (I_n + \lambda M)^{-1} P \quad (22)$$

as the estimated base mesh of the true underlying surface  $S$ .

## Choice of the smoothness parameter $\lambda$

We adopt the *generalized cross validation* (GCV) to determine the smoothness parameter  $\lambda$  in our DLRS model. Specifically, the merit function of GCV is defined as

$$\text{GCV}_n(\lambda) = \frac{\frac{1}{n} \|P - \hat{S}_n(\lambda)\|_F^2}{(1 - \frac{1}{n} \text{tr}[A_n(\lambda)])^2}, \quad (23)$$

where  $A_n(\lambda) = (I_n + \lambda M)^{-1}$ . The optimal value of  $\lambda$  can be computed by minimizing the above GCV function, i.e.,

$$\hat{\lambda}_G = \arg \min_{\lambda > 0} \text{GCV}_n(\lambda). \quad (24)$$

## Asymptotic properties

Denote the error between the estimated base mesh surface and the true underlying surface as

$$r_n(\lambda) = \frac{1}{n} \|\hat{S}_n(\lambda) - S\|_F^2.$$

**Theorem 1.** Assume that  $P$  is the equidistributed sample of a  $C^2$ -smooth surface  $S$ . As  $n \rightarrow \infty$  and  $\lambda \sim n^{-2/3}$  is chosen, we have

$$\mathbb{E}[r_n(\lambda)] = O(n^{-\frac{2}{3}}).$$

# Asymptotic properties

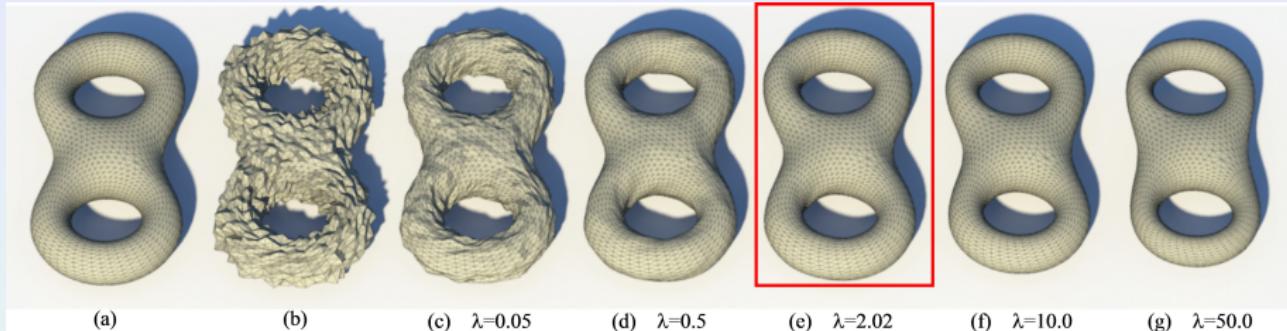
**Theorem 2.** If the smoothness parameter  $\hat{\lambda}_G$  is the GCV choice according to (24), then the estimated base surface  $\hat{S}_n(\hat{\lambda}_G)$  from our DLRS model is asymptotically optimal, i.e.,

$$\frac{r_n(\hat{\lambda}_G)}{\inf_{\lambda \in \mathbb{R}_+} r_n(\lambda)} \rightarrow_p 1,$$

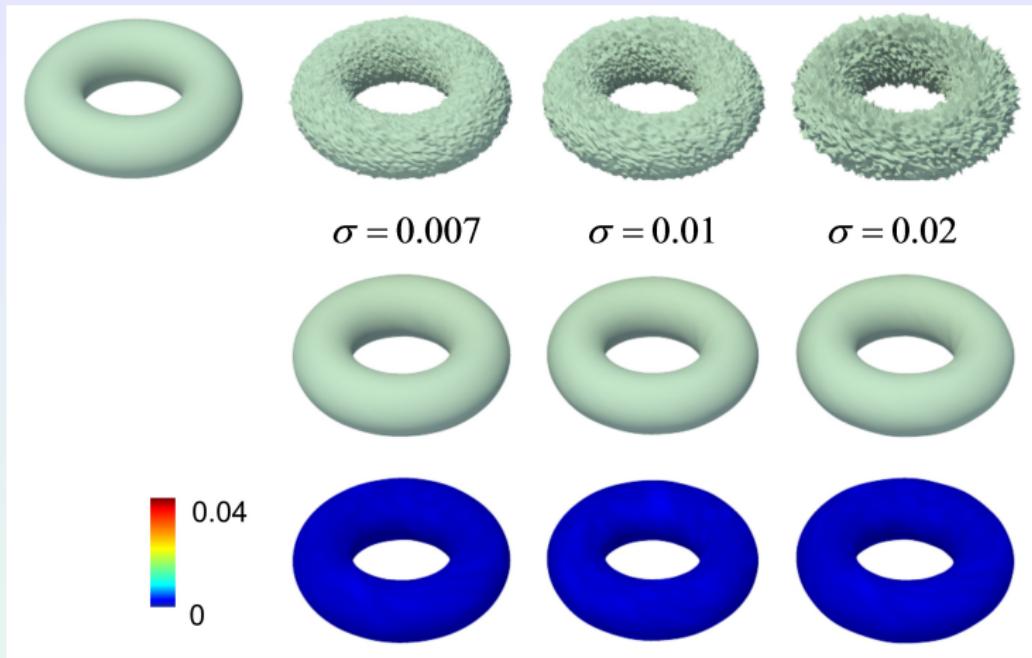
where  $\rightarrow_p$  means the convergence in probability.

## Asymptotic properties

From the above theorems, it is seen that the estimated base surface  $\hat{S}$  asymptotically converges to the ground truth surface  $S$  with probability one as the sample size goes to infinity.



**Figure:** Our approach computes the optimal smoothness parameter to denoise the mesh. (a) the ground truth 8-like model which is a  $C^2$ -smooth surface; (b) the model artificially corrupted by severe synthetic noise; (c)-(g) denoised results by the global smoothing approach with various parameters  $\lambda$  in (21). Our approach obtains the best smoothing result (e) using the optimal parameter  $\lambda = 2.02$ .



**Figure:** Our approach is robust to different amount of noises. Example: a  $C^2$ -smooth torus shape. The first row shows the ground truth model followed by the noisy models with different amount of noise (the variance is shown below each model), the second row shows the denoised models, and the third row visualizes the differences between the denoised models and the ground truth model.

# Residual

By the denoising phase we now have the estimated base surface  $\hat{S} = \{\hat{s}_i\}_1^n$  from the input noisy mesh  $P = \{p_i\}_1^n$ . The residual between  $\hat{S}$  and  $P$  is defined as

$$b_i = (p_i - \hat{s}_i)^T \hat{n}_i, \quad i = 1, \dots, n, \quad (25)$$

where  $\hat{n}_i$  is the unit normal vector of surface  $\hat{S}$  at  $\hat{s}_i$ . Denote  $b = (b_1, \dots, b_n)^T$  as the residual vector.

## Residual

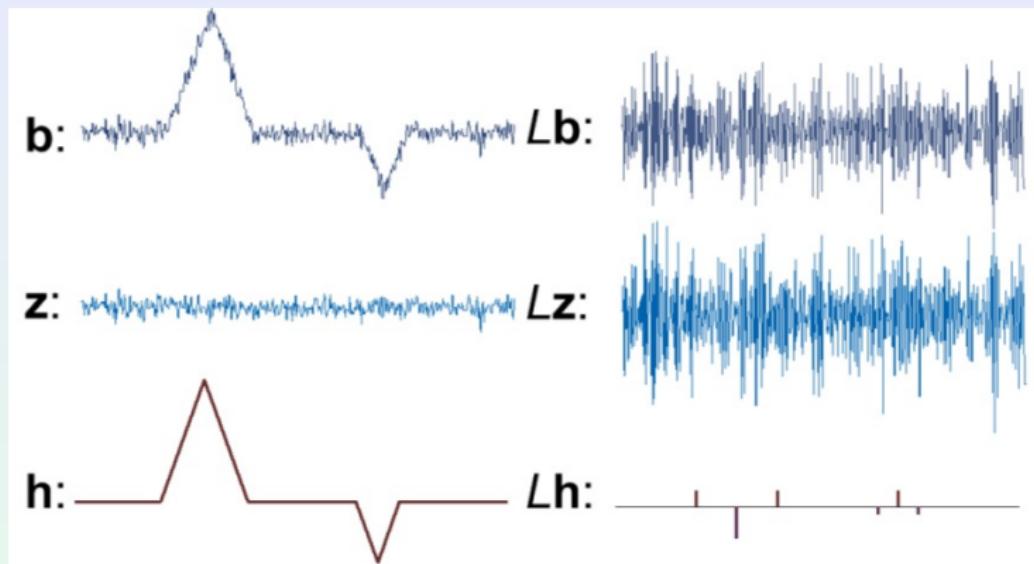
As an inference of asymptotic properties, the residual signal  $b$  is essentially i.i.d. noise when the true underlying surface of input mesh  $P$  is at least  $C^2$ -continuous as a whole.

But in case the underlying surface containing sharp features, the residual  $b$  inevitably mixes the information from features and the noise, as illustrated in Figure 4 (left), and is decomposed as

$$b = h + z, \quad (26)$$

where  $h = (h_1, \dots, h_n)^T$  is the unknown signal of the features and  $z = (z_1, \dots, z_n)^T$  is the measurement errors.

# Residual



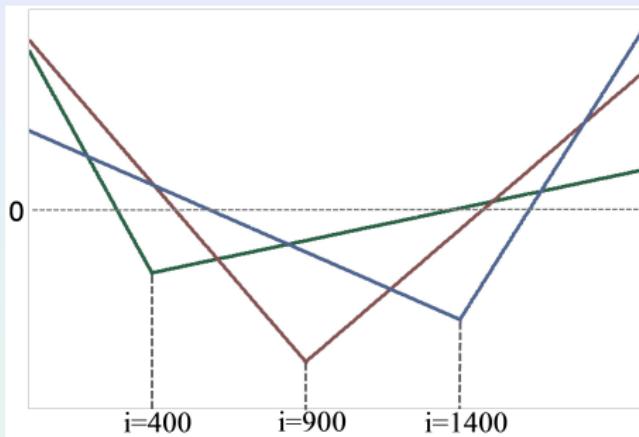
**Figure:** 1D illustration of signals and their Laplacians. Left: the residual  $b$  (upper) is a mixture of noises  $z$  (middle) and features  $h$  (lower). Right: the corresponding Laplacians of signals on the left.

## Coherent dictionary for shape features

Considering the signal of shape feature  $h$  in our case, we discover that shape features can be represented as sparse in some coherent dictionary which is constructed by the pseudo-inverse matrix  $L^+$  of the Laplacian matrix  $L$  of the shape.

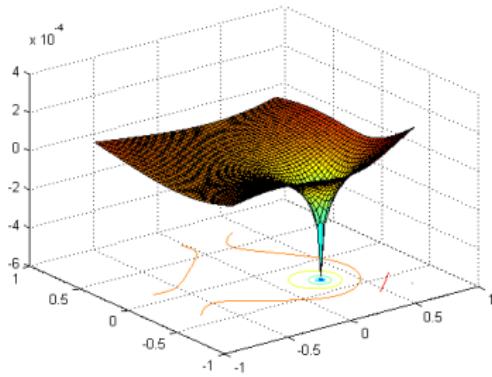
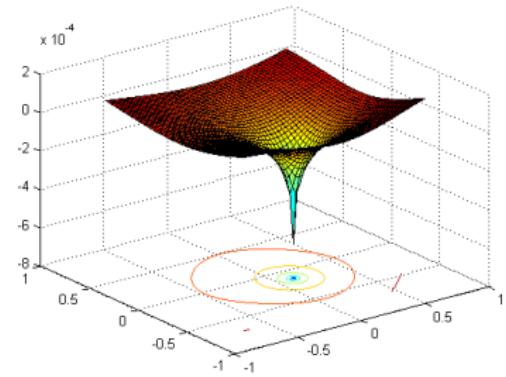
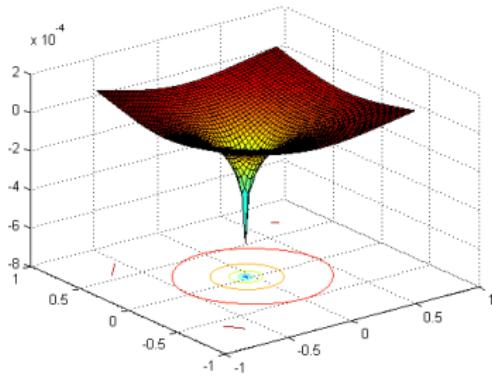
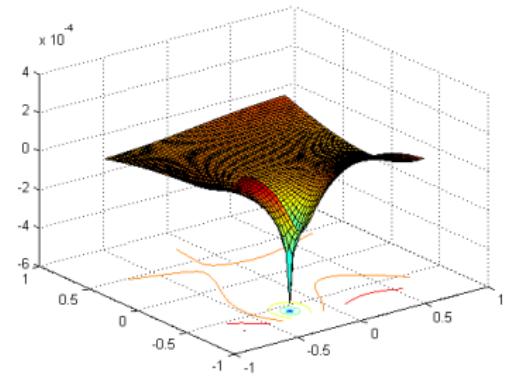
It is not surprising to see that the Laplacian of feature  $Lh$  is indeed a sparse signal and has quickly decaying coefficients, as illustrated in Figure 4 (right).

# Coherent dictionary for shape features



**Figure:** Three basis functions corresponding to the 400-th, 900-th, and 1400-th columns of  $L^+$  respectively. This evidence suggests that  $L^+$  is a coherent dictionary for representing the feature signal on the shape. We use 1D case as an illustration.

# Coherent dictionary for shape features



## $\ell_1$ -analysis on residual

As we have found a coherent dictionary  $D = L^+$  for representing  $h$  sparsely, i.e.,  $D^+h$  is sparse, we can thus formulate the problem of recovering feature signals  $h$  from the residual  $b$  as an  $\ell_1$ -analysis compressed sensing optimization.

In our setting, the sensing matrix is the identity matrix, and the coherent dictionary is  $D = L^+$ . Thus we have the following  $\ell_1$ -analysis compressed sensing optimization:

$$\min_h \|Lh\|_1 \quad \text{s.t.} \quad \|h - b\|_2 \leq \epsilon, \quad (27)$$

where  $\epsilon^2$  is a likely upper bound on the noise power  $\|z\|_2^2$ .

## $\ell_1$ -analysis on residual

The roles of the penalty and the constraint in (27) might also be reversed if we choose to constrain the sparsity and obtain the best fit for that sparsity. Here we prefer to solve an equivalent optimization:

$$\min_{\mathbf{h}} \|\mathbf{h} - \mathbf{b}\|_2^2 \quad \text{s.t.} \quad \|L\mathbf{h}\|_1 \leq \tau, \quad (28)$$

where  $\tau$  is a tune parameter controlling the sparsity. We call  $\tau$  the sparsity parameter.

## Weighted $\ell_1$ -analysis

Consider the weighted  $\ell_1$ -analysis on the residual

$$\begin{aligned} \min_{\mathbf{h}} \quad & \|\mathbf{h} - \mathbf{b}\|_2^2 \\ \text{s.t.} \quad & \|W(L\mathbf{h})\|_1 = \sum_{i=1}^n w_i |L_i \mathbf{h}| \leq \tau \end{aligned} \tag{29}$$

where  $W = \text{diag}(w_1, \dots, w_n)$  and  $w_1, \dots, w_n$  are positive weights. The weighted  $\ell_1$ -analysis optimization (29) can be regarded as a relaxation of an  $\ell_0$ -minimization problem.

## Weighted $\ell_1$ -analysis

It is desired that the weights could be to counteract the influence of the signal magnitude on the  $\ell_1$ -penalty function. Ideally, the weights are inversely proportional to the true signal magnitude, i.e.,

$$w_i = \begin{cases} \frac{1}{|L_i h|} & L_i h \neq 0, \\ \infty & L_i h = 0. \end{cases} \quad (30)$$

The large entries in  $W$  force the Laplacian of feature  $Lh$  to concentrate on the indices where  $w_i$  is small. These constructed weights precisely correspond to the indices where  $Lh$  is nonzero.

It is of course impossible to construct the precise weights (30) without knowing the feature signal  $h$  itself, but this suggests more generally that large weights could be used to discourage nonzero entries in the recovered  $Lh$ , while small weights could be used to encourage nonzero entries.

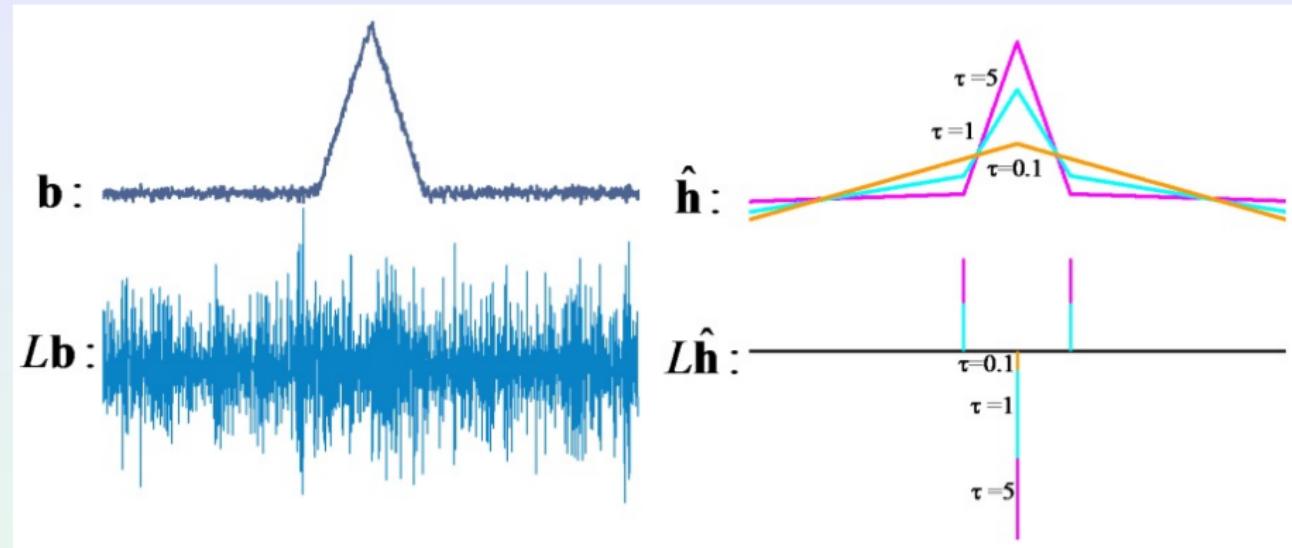
## Weighted $\ell_1$ -analysis

Based on the geometric information of the estimated base surface  $\hat{S}$ , we design the weights as follows

$$w_i = \frac{1}{\rho + \|L_i \hat{S}\|}, i = 1, \dots, n,$$

where  $\rho$  is a small number ( $\rho = 10^{-7}$  by default) that provides numerical stability and should be set slightly smaller than the expected nonzero magnitudes of  $Lh$ . With these well designed weights, we then perform a weighted  $\ell_1$ -analysis (29) on the residual for recovering the feature signal.

# Weighted $\ell_1$ -analysis



**Figure:** 1D illustration of results using different sparsity parameters  $\tau$  in the  $\ell_1$ -analysis optimization. Left: the residual signal  $\mathbf{b}$  (upper) and its Laplacian (lower); Right: the recovered feature signals  $\hat{\mathbf{h}}$  corresponding to different values 0.1 (orange), 1 (blue), and 5 (magenta) of  $\tau$  (upper) and their Laplacians (lower). It is seen that only prominent feature is identified using small  $\tau = 0.1$  and more features are identified using larger  $\tau = 1$  and  $\tau = 5$ . The  $\ell_1$ -analysis

## Iterative feature recovering

If we set large value of the sparsity parameter  $\tau$ , the optimization (29) will recover features but may introduce some non-feature points in the result. Thus we prefer to choose a small value of  $\tau$ . For some models with large portion of features, the solution to (29) returns only the most prominent (sharpest) features.

To recover the other features, our idea is to modify the rows corresponding to the identified features in the Laplacian matrix  $L$  and perform the weighted  $\ell_1$ -analysis optimization (29) in an iterative manner. In each iteration, we identify a part of features. After marking the identified features, we find more features in the next iteration. This is reasonable as features become sparser in the latter iterations as the identified features in previous iterations have already marked and do not contribute to the feature recovering.

## Feature types and classifications

Generally there are two types of sharp features, corners and creases, as shown in Figure 7, on 3D shapes. A corner point is the one at which the tangent of any passing curve on the surface is discontinuous. A crease curve introduces the discontinuities of first derivatives across it, but preserves  $C^2$ -continuity along it. Corners can also be the intersections of several creases.

After identifying locations of the features by the weighted  $\ell_1$ -analysis (29) on the residual, we adopt a simple scheme to classify their types. If a feature is isolated from others, it is identified as a corner. If a feature point has a few feature points in its neighbor, we compute a dominant direction by PCA and classify it as a crease along this direction, see Figure 7.

# Feature types and classifications

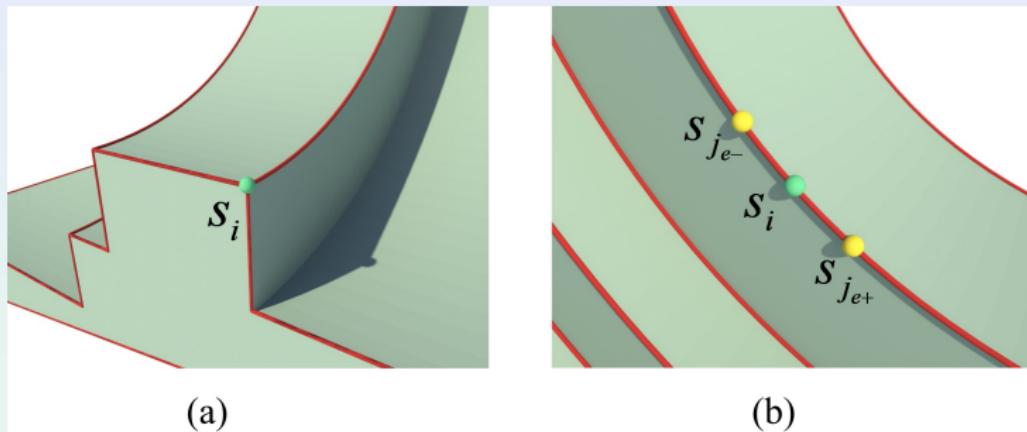


Figure: Two types of feature points: (a) corner; (b) crease.

## Feature aware modification of Laplacian matrix

If the vertex  $s_i$  is identified as a corner, we remove the Laplacian penalty  $\mathcal{L}(s_i) = L_i S$  by setting the  $i$ -th row of Lapalcian matrix  $L_i = 0$ . If the vertex  $s_i$  is identified as a point on a crease  $\mathcal{E}$ , we only remove the straddling smoothness penalties and yield a term of the form

$$\mathcal{L}(s_i) = L_i S = \sum_{j \in \mathcal{N}(i) \cap \mathcal{E}} (s_j - s_i) = s_{j_{e-}} + s_{j_{e+}} - 2s_i,$$

where  $s_{j_{e-}}$  and  $s_{j_{e+}}$  are the adjacent neighbors of  $s_i$  along the crease  $\mathcal{E}$ , as illustrated in Figure 7(b).

## Termination condition of iterations

We adopt a statistical method of nonparametric test for checking whether the residual  $b$  contains more features. Specifically, residuals were randomly divided into two sets  $b_1$  and  $b_2$ . We use a two-sample Kolmogorov-Smirnov test to compare the distributions of the values in the two sets  $b_1$  and  $b_2$ . We state that

- The null hypothesis  $H_0$ :  $b_1$  and  $b_2$  are from the same continuous distribution.
- The alternative hypothesis  $H_1$ : they are from different continuous distributions.

This hypothesis does not specify what that common distribution is (e.g. normal or not normal). The result is 1 if the test rejects the null hypothesis at the  $\alpha$  significance level; 0 otherwise. We use the significance level  $\alpha = 0.05$  in our implementation.

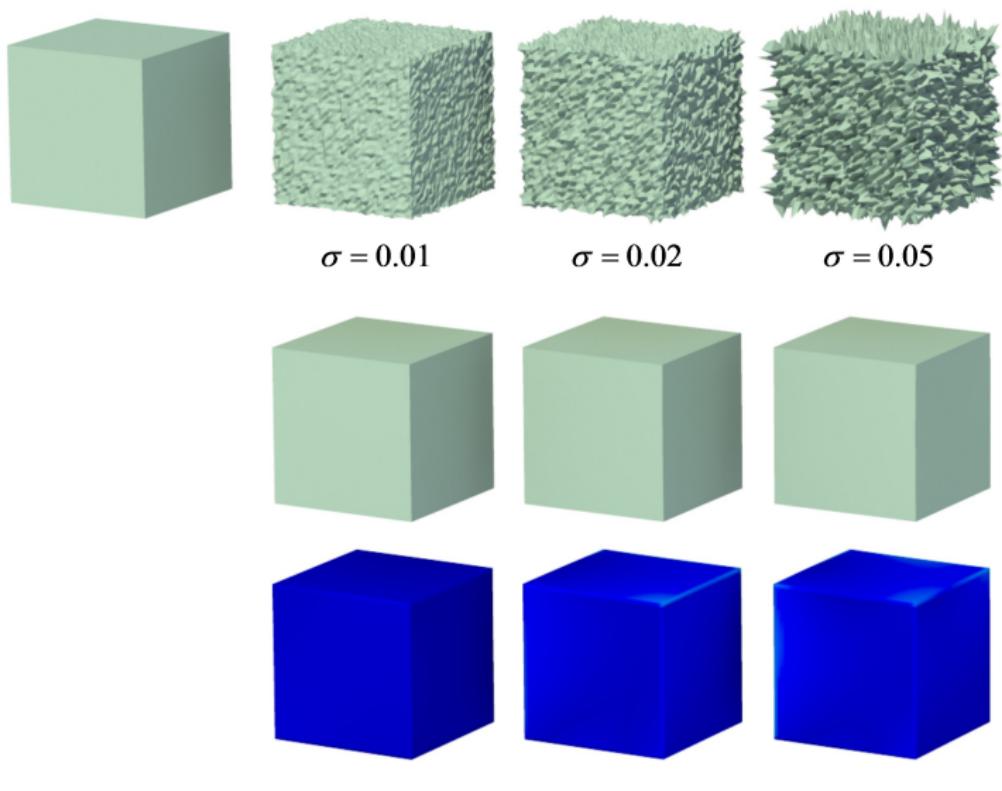
## Pseudo-code

The pseudo-code of the iterative feature recovering via  $\ell_1$ -analysis can be found in the following algorithm.

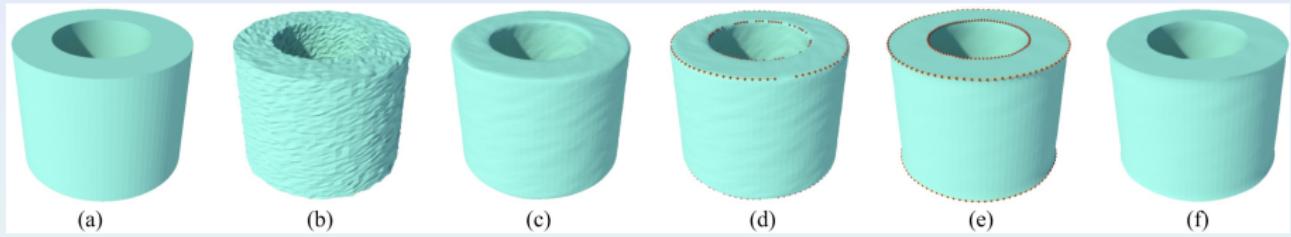
**Input:** the points  $P$ , its Laplacian  $L$ , and the sparsity parameter  $\tau$

**Output:** denoised mesh with features

- ① Call  $\hat{S} = \text{DLRS}(P, L)$ .
- ② Calculate the residual  $b$  according to (25).
- ③ If the result of the two-sample Kolmogorov-Smirnov test on residual  $b$  is 1, go to Step 4; Otherwise, exit and output the current  $\hat{S}$ .
- ④ Recover the features  $\hat{h}$  by the weighted  $\ell_1$ -analysis (29), and get the reliable locations of  $\hat{h}$  indicated by  $L\hat{h}$ .
- ⑤ Classify the features  $\hat{h}$  and modify accordingly the Lapalcian matrix  $L$  based on their feature types.
- ⑥ Go to Step 1.



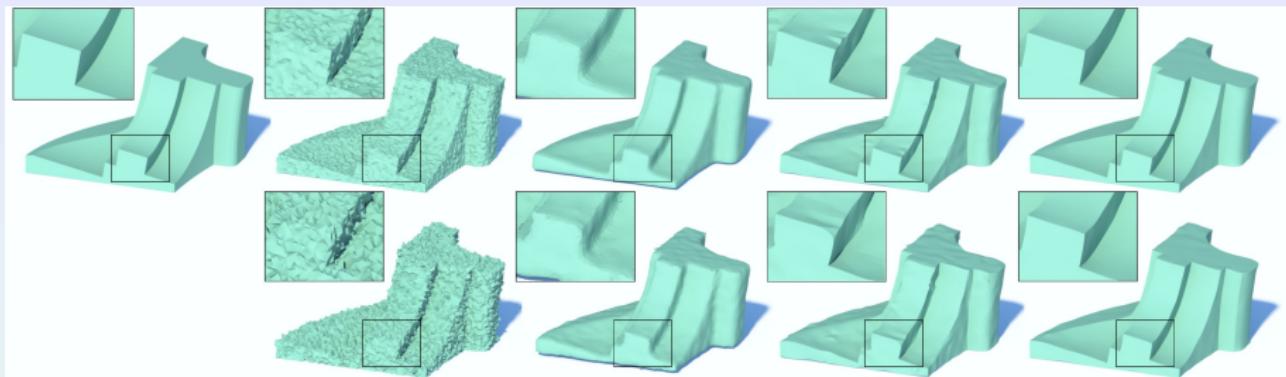
**Figure:** Our approach is robust to different amount of noises. Example: a cube model with sharp features.



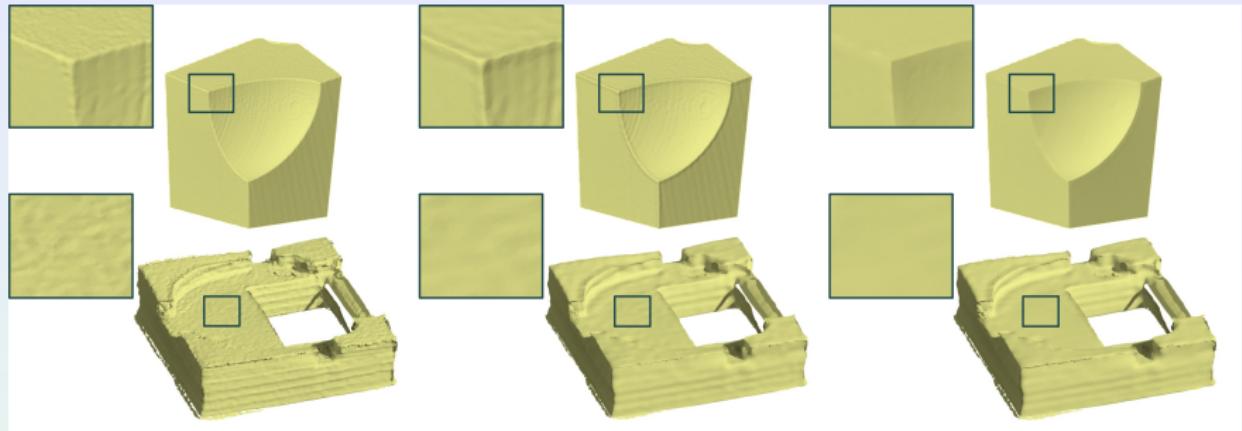
**Figure:** Our algorithm progressively recovers sharp features. (a) the ground truth model; (b) the noisy model; (c) the denoised model using our approach; (d)-(e) progressively recovered features shown in red dots with two iterations; (f) the final result.



**Figure:** Main phases of our approach. (a) The input noisy mesh model; (b) phase I: denoising the model using the DLRS scheme; (c) phase II: recovering sharp features progressively based on the  $\ell_1$ -analysis optimization; (d) the final denoised result with recovered features.



**Figure:** Comparisons with previous approaches. From left to right: the ground truth Fandisk model, the models corrupted by different levels of artificial noise ( $\sigma = 0.005$  in the upper row and  $\sigma = 0.01$  in the lower row), denoised results with vertex based bilateral mesh filtering [Fleishman et al. 2003], normal based bilateral mesh filtering [Zheng et al. 2010], and our approach. The small regions with the frames are magnified to clearly show the differences. All meshes are flat-shaded to show faceting.



**Figure:** Applying our approach to scanning data of real objects. From left to right: the input scanning raw data; the denoised results using the non-iterative smoothing method [Jones et al. 2003]; the denoised results using our approach. The close-up views show the details in framed regions.

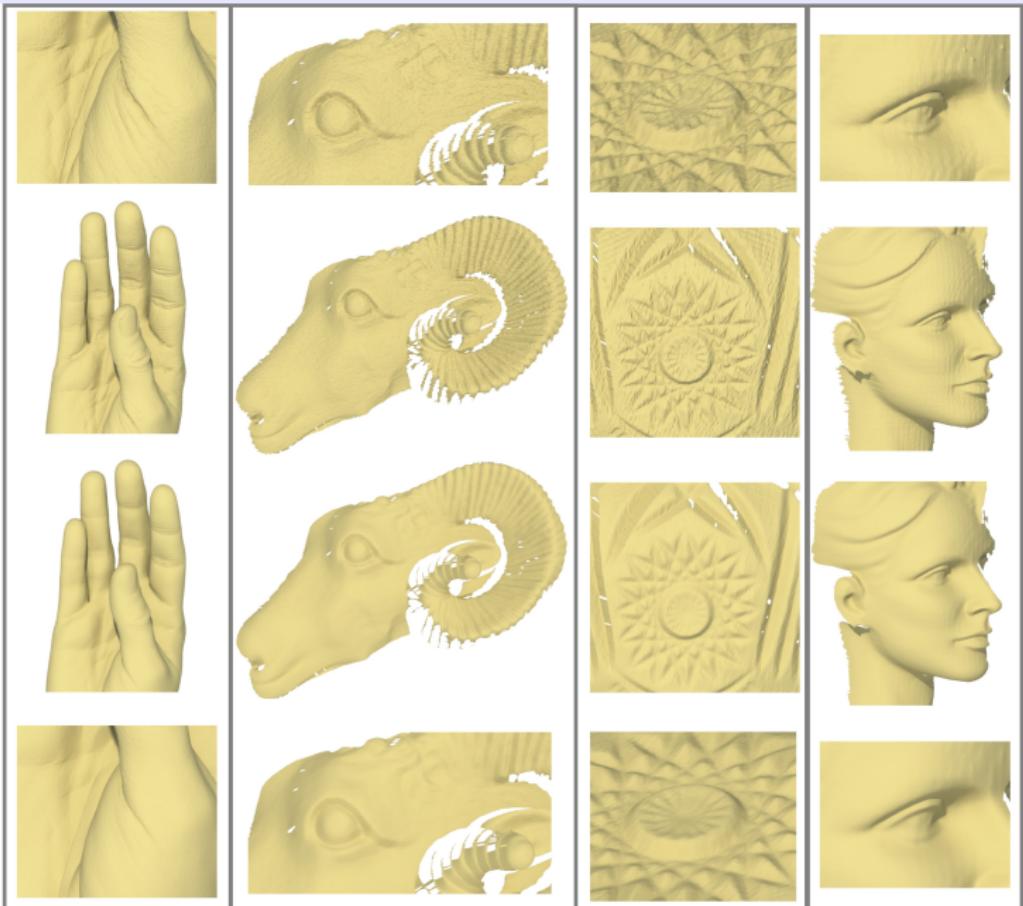


Figure: Applying our approach to real point clouds. The second row is the

yangzw@ustc.edu.cn (USTC)

Sparsity

# Summary

A two-phase approach for decoupling features and noises on discrete surfaces:

- The first phase generates a base mesh which is obtained by denoising the input data using the DLRS estimator.
- The second phase recovers sharp features from the residual between the base mesh and the input mesh based on an  $\ell_1$ -analysis compressed sensing optimization.

# Summary

If the true underlying surface of input data is  $C^2$  smooth, Step 1 obtains a smoothed data which is theoretically guaranteed to asymptotically converge to the true underlying surface (Theorem 2).

When the underlying surface is not a  $C^2$ -smooth surface, e.g., a piecewise  $C^2$  surface with  $C^0$  features, Step 2 is performed to extract features from residual between original and smoothed data. By removing the unnecessary smoothness penalty at locations of features, Step 1 performs equivalently as denoising each  $C^2$  patch individually and thus eventually obtains the optimal approximation to the underlying surface with theoretical guaranty.

# Contributions

- **Asymptotically optimal surface denoising:** Our denoising approach is fully automatic without tuning any parameter. The smoothness parameter is automatically computed. The denoised surface is guaranteed to asymptotically converge to the true underlying surface with probability one if the underlying surface is  $C^2$ -smooth.
- **Faithful feature recovering by  $\ell_1$ -analysis:** We successfully apply the  $\ell_1$ -analysis compressed sensing technique to identify and recover sharp features from the residual between the estimated base surface and the input surface. This is based on our discovery that the pseudo-inverse matrix of the Laplacian matrix of the mesh is a coherent dictionary for sparsely representing the sharp feature signal on the shape.

# Outline I

- ① Compressed Sensing
- ② Sparse Modeling
  - Sparsity-seeking representations
  - Numerical optimization
  - Applications
- ③ Sparse Optimization
  - Sparse Optimization Models
  - Sparse Optimization Algorithms
- ④ Decoupling Noises and Features via  $\ell_1$ -analysis Compressed Sensing
  - Discrete Laplacian regularization smoothing
  - Feature recovering via  $\ell_1$ -analysis optimization

# Outline II

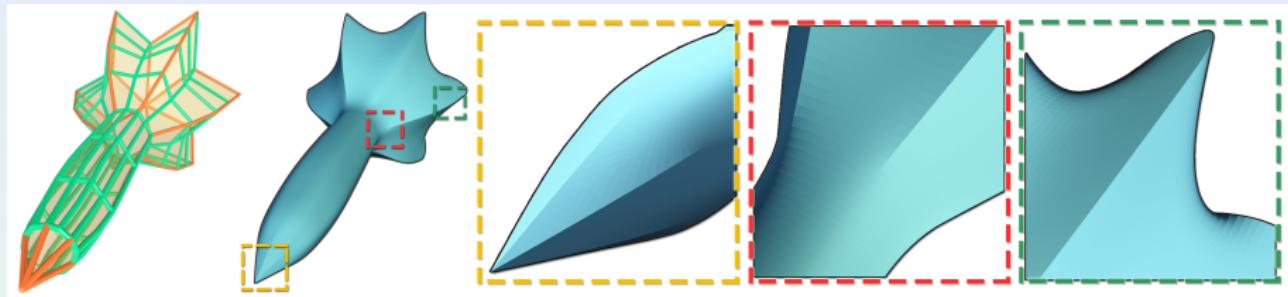
- Experimental results

5 Construction of Manifolds via Consistent Sparse Representations

6 Sparse Representation with Parameterization Optimization

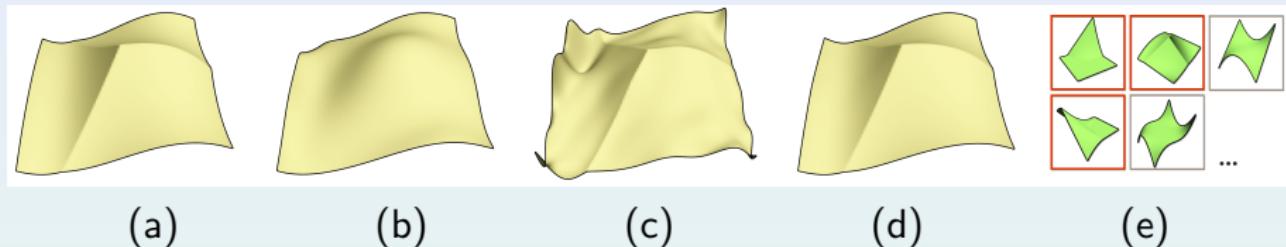
7 Conclusion and Future Work

# Motivation



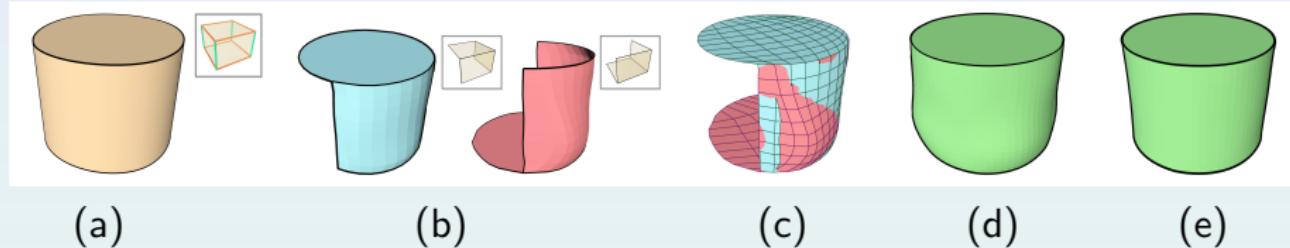
**Figure:** To construct a manifold that approximates various features like crease, dart, and cusp well.

# Motivation



**Figure:** Given a set of sampling points from a  $C^0$  surface (a), the obtained fitting surface via polynomials (b), redundant atom functions (c), and sparse representation (d). (e) The first five atom functions ( $C^0$  shape functions in red boxes and polynomials in gray boxes) that are selected in the sparse representation.

# Motivation



**Figure:** Given subdivision surface with control mesh (a), the results of sparse representations on individual charts (b), the surface's appearance in overlapping regions (c), the result of a manifold (d) by simple blending, and the obtained surface (e) by consistent sparse representation.

# Manifold

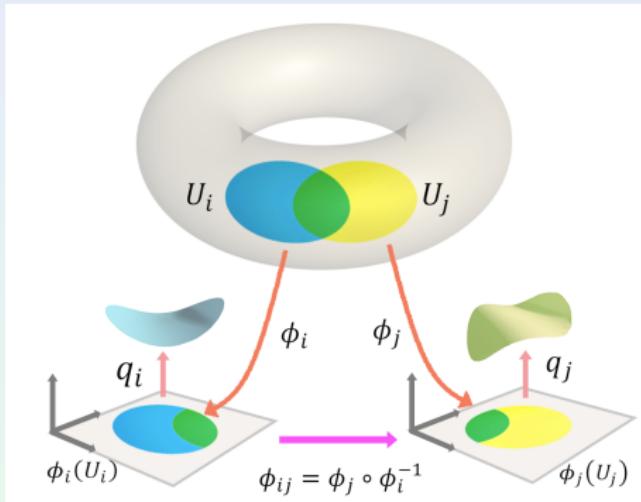


Figure: The definition of a manifold.

# The $C^0$ atom functions

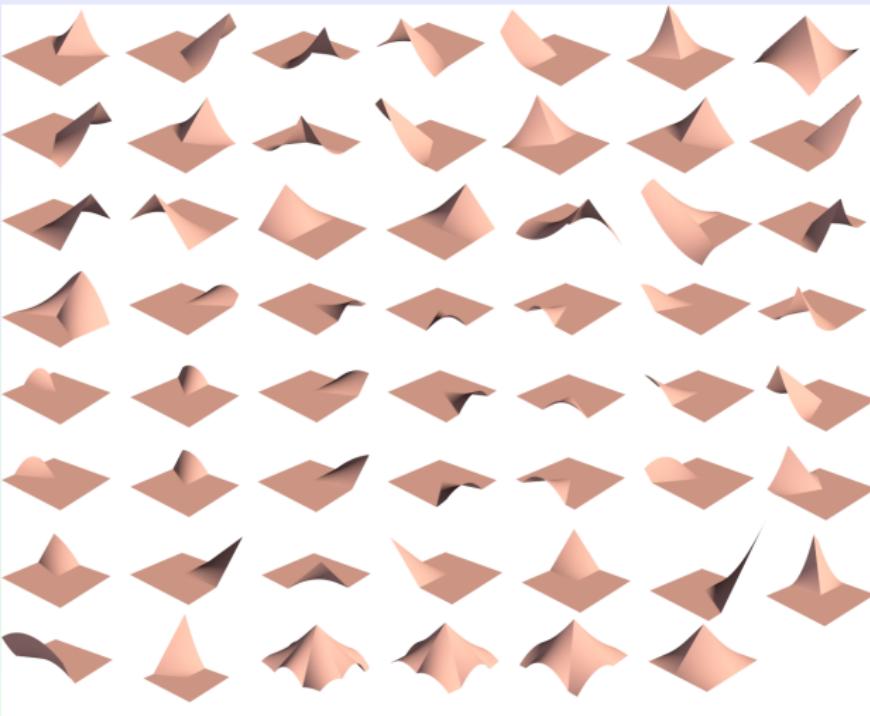


Figure: 55 shape functions with valence of 3 to 7 (which is suitable for most domain manifold).

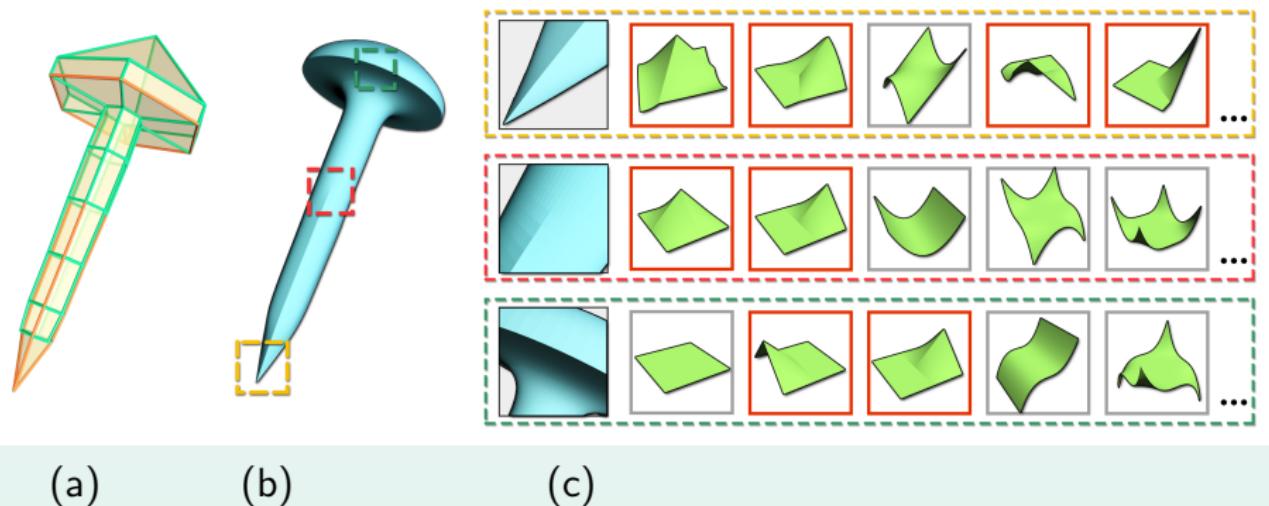
# Compatible sparse representations

$$E_{\text{fit}}(C) = \sum_{k=1}^K \left( \sum_{j: p_k \in U_j} \left( w_j(u_k^j, v_k^j) \sum_{l=1}^L c_{lj} a_l(u_k^j, v_k^j) \right) - h_k \right)^2 \quad (31)$$

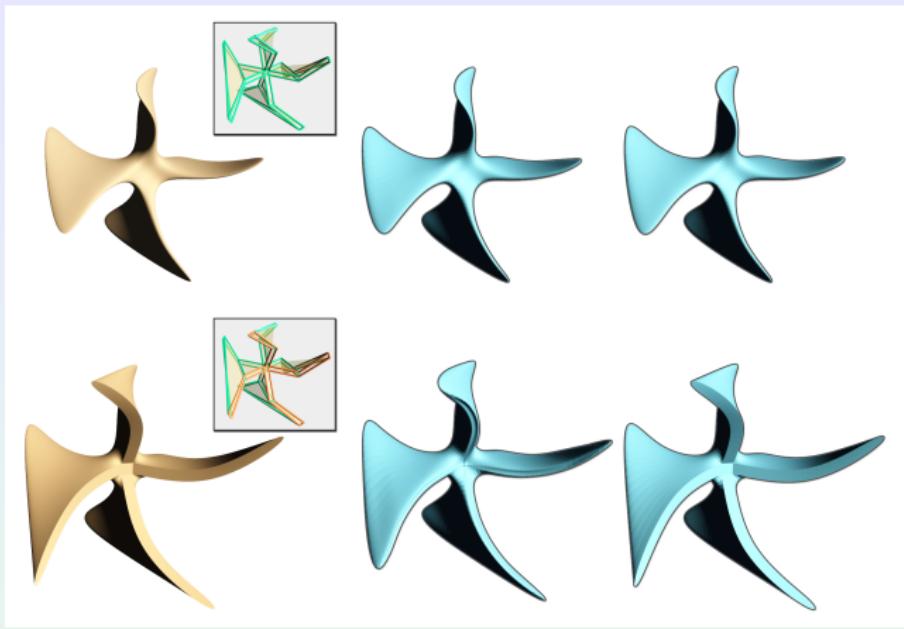
where  $c_j = (c_{1j}, \dots, c_{Lj})^T$  is the coefficient vector of the local function on the chart  $(U_j, \phi_j)$  and  $C = (c_1, \dots, c_n)$  is the matrix of decision variables.

To obtain the sparse coefficients of the local function on each chart, we formulate the sparse representation of global geometry function as

$$\begin{aligned} & \min_C E_{\text{fit}}(C) \\ & \text{s.t. } \|c_j\|_0 \leq \delta, j = 1, 2, \dots, n. \end{aligned} \quad (32)$$



**Figure:** Constructing manifold to approximate subdivision surface. (a) The input control mesh with enforced  $C^0$  features shown in red; (b) The manifold generated by our approach; (c) The close-up views of three parts of sharp features including crease, dart, and cusp. The first five atom functions that are adaptively selected for representing the local features are shown.

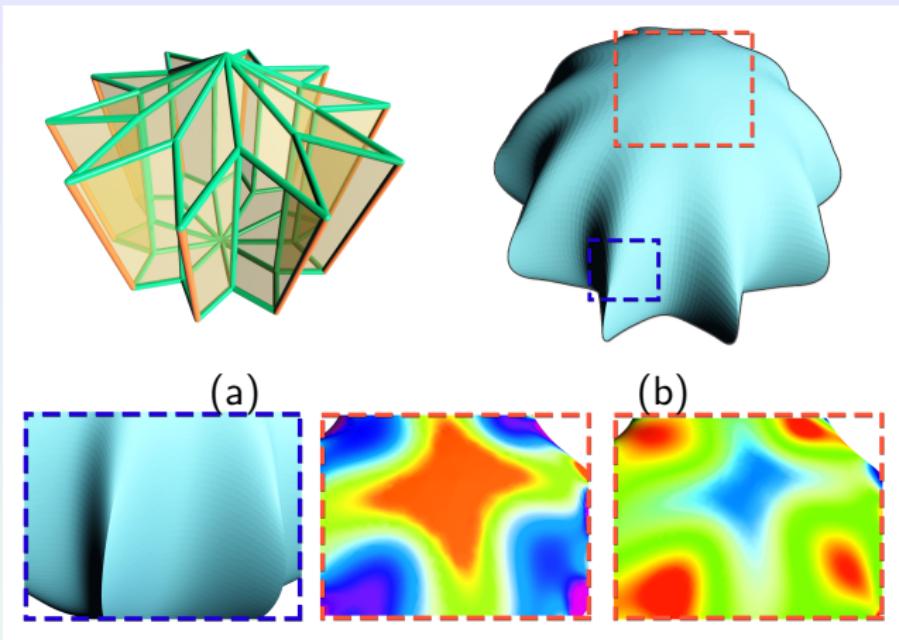


(a)

(b)

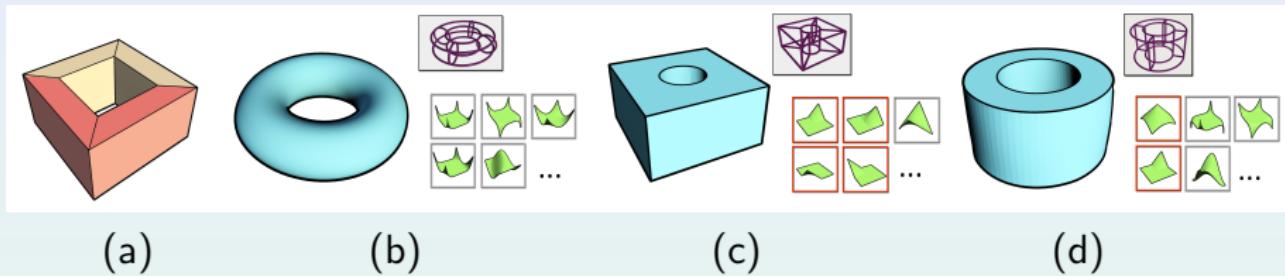
(c)

**Figure:** (a) The limit subdivision surfaces of input control meshes (in the upper right) where edges in red are enforced with  $C^0$  constraints; (b) manifolds produced by the method of [Ying and Zorin 2004]; (c) manifolds generated by our approach. Upper row: the results from smooth subdivision surface; Lower row: the results from subdivision surface with sharp features.



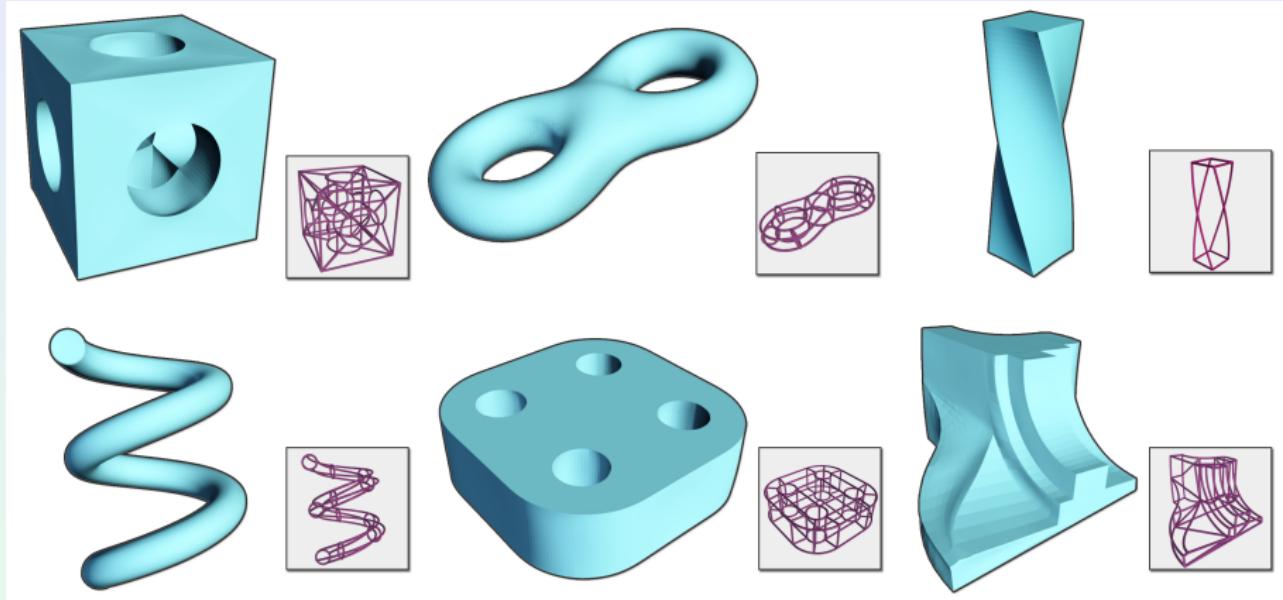
**Figure:** (a) A control mesh with  $C^0$  constraints shown in red. The top vertex of (a) has a valence of 10. We discard the  $C^0$  atom functions from the dictionary at the chart around this vertex. The manifold surface constructed by our approach is shown in (b). The surface adapts to the sharp features in blue region and is  $C^\infty$  in the red region. (c) the close-up view of blue region; (d,e) Gaussian curvature and mean curvature of the red region.

# Application to curve network



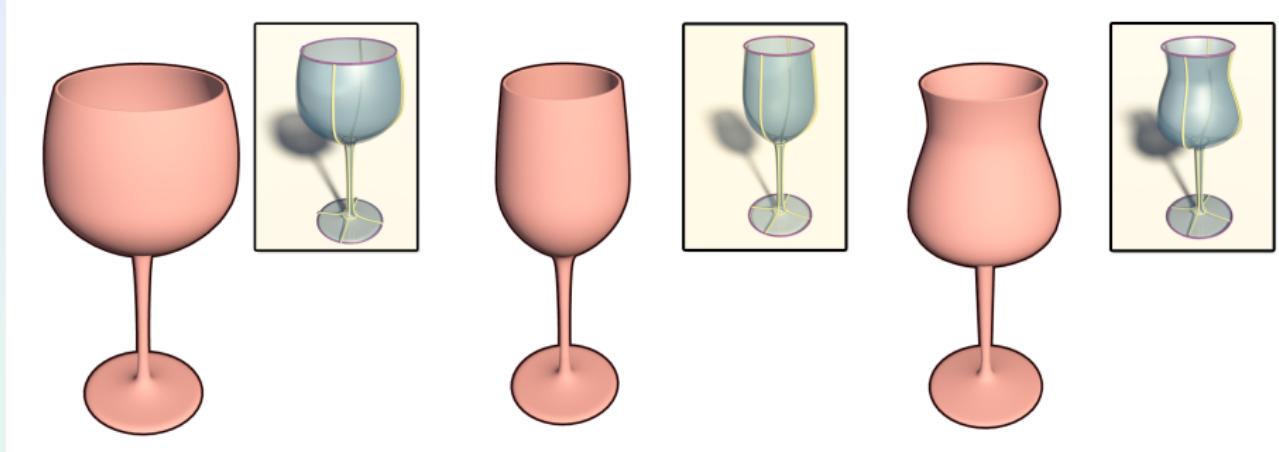
**Figure:** Generating manifold from different geometries defined on the same domain manifold. (a) The domain mesh; (b-d) three manifolds created from given curve networks (in the upper right) and the first five chosen atom functions (in the lower right).  $C^0$  shape functions are shown with red frames.

# Application to curve network



**Figure:** Examples of manifold surfaces generated from 3D curve networks (in the lower right).

# Application to curve network



**Figure:** Three different wine glasses obtained by our method with specifying smoothness across some edges (yellow ones).

# Outline I

- ① Compressed Sensing
- ② Sparse Modeling
  - Sparsity-seeking representations
  - Numerical optimization
  - Applications
- ③ Sparse Optimization
  - Sparse Optimization Models
  - Sparse Optimization Algorithms
- ④ Decoupling Noises and Features via  $\ell_1$ -analysis Compressed Sensing
  - Discrete Laplacian regularization smoothing
  - Feature recovering via  $\ell_1$ -analysis optimization

# Outline II

- Experimental results

5 Construction of Manifolds via Consistent Sparse Representations

6 Sparse Representation with Parameterization Optimization

7 Conclusion and Future Work

# Motivation & Purpose

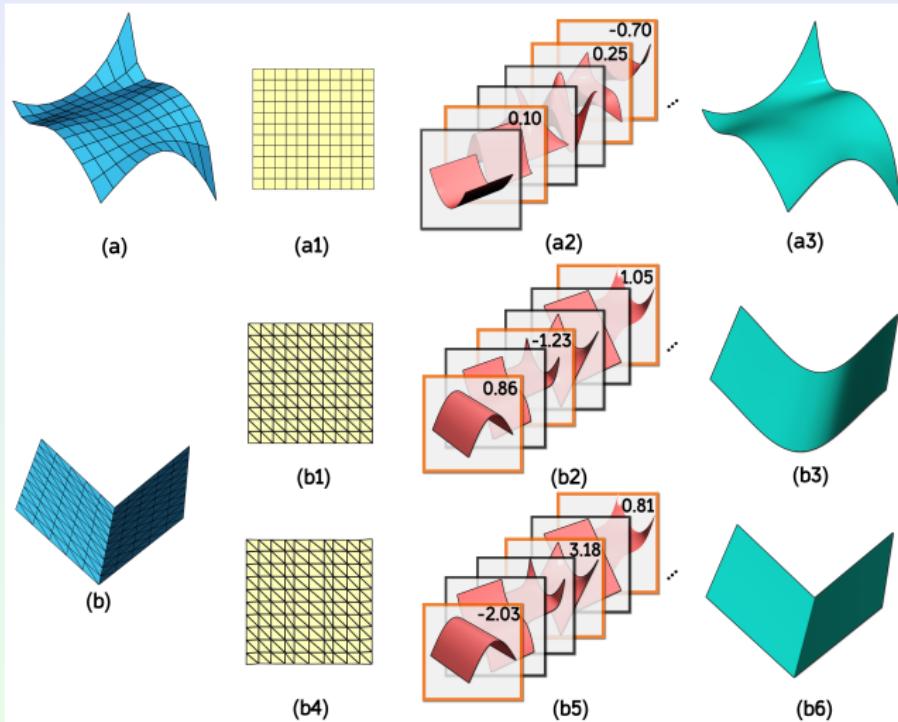
- Fitting geometry data based on a sparse optimization can efficiently overcome the overfitting artifacts.
- The parameterization of input data is simultaneously optimized to respect sharp features of target geometry.

# Composite function

$$f(u) = h \circ \Phi(u), \quad u \in M \quad (33)$$

$$\begin{cases} \Phi : M \rightarrow M, \\ h : M \rightarrow \mathbb{R}. \end{cases}$$

# Composite function



## Fitting problem

Denote  $\mathbf{b} = (b_1, \dots, b_m)^T$  as a set of scalar observation from some underlying Euclidean space  $\mathbb{R}^k$ .

Suppose  $D = \{d_j(u)\}_{j=1}^n$  is a family of  $n$  basis functions defined on  $\mathbb{R}^k$ .

The fitting problem is to find the coefficients  $\mathbf{c} = (c_1, \dots, c_m)^T$  by minimizing the error function

$$F(\mathbf{c}, U) = \sum_{i=1}^m \left[ b_i - \sum_{j=1}^n c_j d_j(u_i) \right]^2,$$

where  $U = (u_1, \dots, u_m) \in \mathbb{R}^{k \times m}$  are the underlying parameters of data.

## Sparse fitting

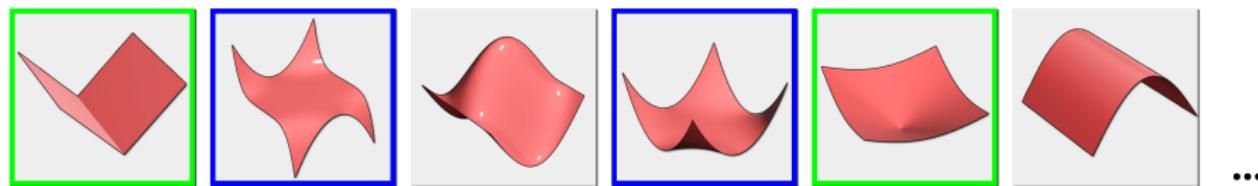
We usually adopt a set of redundant basis functions  $\{d_j(u)\}_{j=1}^n$  and fit the input data by a sparse combination of these functions, i.e.,

$$\begin{aligned} \min_{\mathbf{c}} \quad & F(\mathbf{c}, U) \\ \text{s.t.} \quad & \|\mathbf{c}\|_0 \leq s \end{aligned} \tag{34}$$

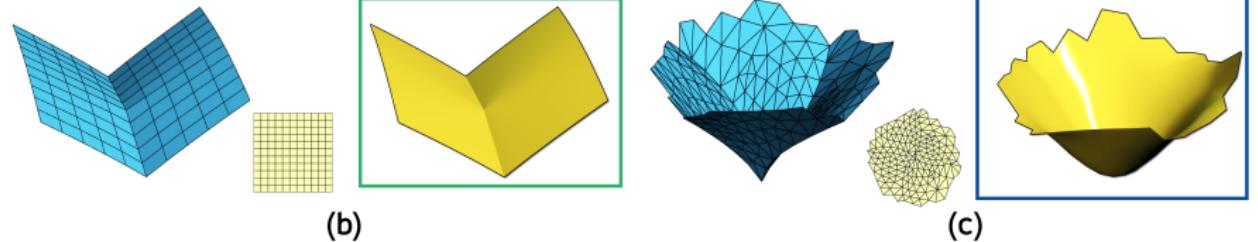
where  $s$  is the sparsity parameter.

With sparse fitting, we can use redundant basis functions and then rely on the sparsity regularization to choose the appropriate basis functions to represent the input geometry data.

# Sparse fitting



(a)



(b)

(c)

# Sparse fitting with parameterization optimization

We optimize the coefficients  $c$  and the parameterization  $U$  of the input data simultaneously as follows:

$$\begin{aligned} \min_{c, U} \quad & F(c, U) \\ \text{s.t.} \quad & \|c\|_0 \leq s. \end{aligned} \tag{35}$$

## Parametric transformation

The parametric transformation is defined by a mapping  $\Phi : \mathbb{R}^2 \rightarrow \mathbb{R}^2$ .

Assume that  $\{u_i\}_{i=1}^m$  is a set of vertices in domain  $M$  and  $\{v_i\}_{i=1}^m$  is a set of vertices in domain  $N$  with  $\Phi(u_i) = v_i$ . Let  $T_M$  be the triangulation of  $\{u_i\}_{i=1}^m$  and  $T_N$  the triangulation of  $\{v_i\}_{i=1}^m$ .

Thus we express the mapping  $\Phi$  as a piecewise linear transformation of two-dimensional triangulated structure, i.e.,

$$\Phi(u) = w_1\Phi(u_{i_1}) + w_2\Phi(u_{i_2}) + w_3\Phi(u_{i_3})$$

for any  $u = w_1u_{i_1} + w_2u_{i_2} + w_3u_{i_3}$  lying in the triangle  $\Delta u_{i_1}u_{i_2}u_{i_3}$  with its barycenter coordinates  $(w_1, w_2, w_3)$ .

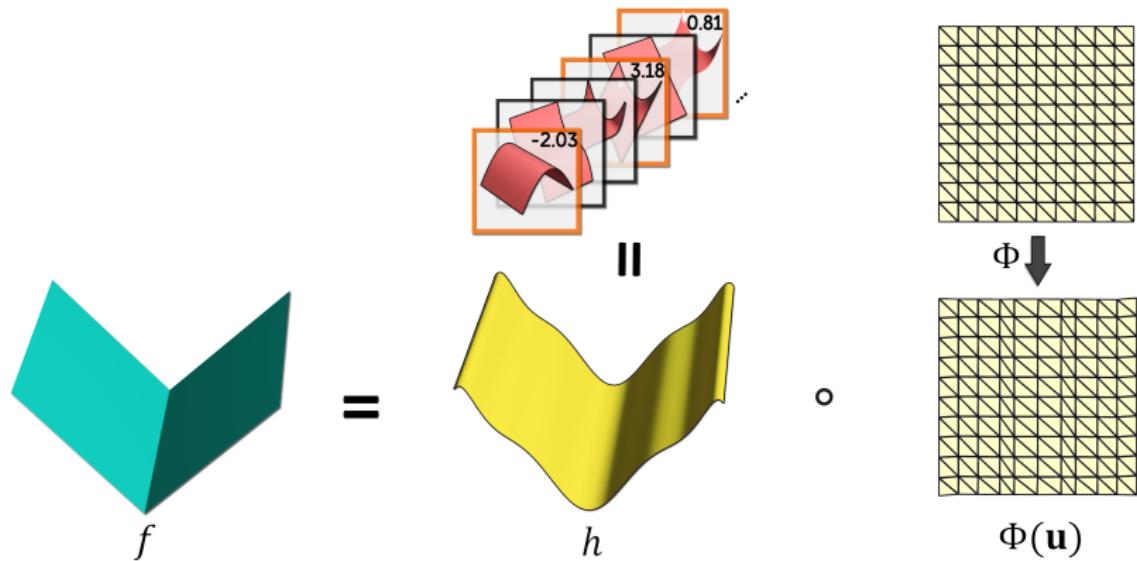
## Composite representation

A function  $f$  defined on domain  $M \subset \mathbb{R}^2$  composites of some function  $h : \mathbb{R}^2 \rightarrow \mathbb{R}$  and a mapping  $\Phi : \mathbb{R}^2 \rightarrow \mathbb{R}^2$ .

Here  $h$  is a linear combination of a certain set of basis functions with  $h(u) = \sum_{j=1}^n c_j d_j(u)$ . The composite function is written as:

$$f(u) = h \circ \Phi(u) = \sum_{j=1}^n c_j d_j(\Phi(u)). \quad (36)$$

# Composite representation



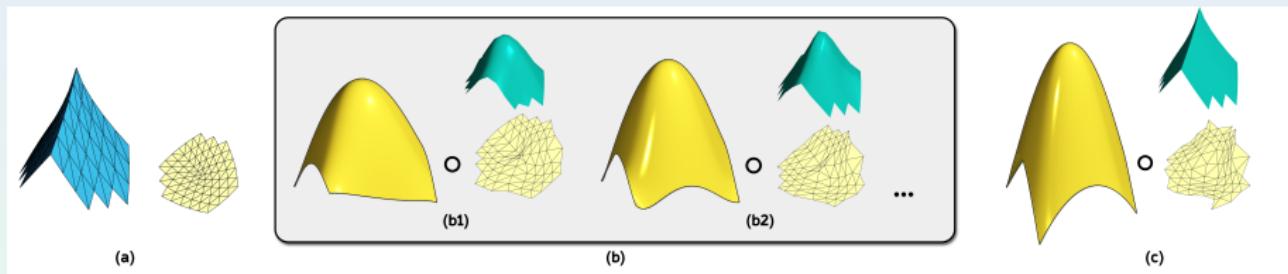
# Optimization framework

We formulate the sparse representation of geometry data as the following optimization:

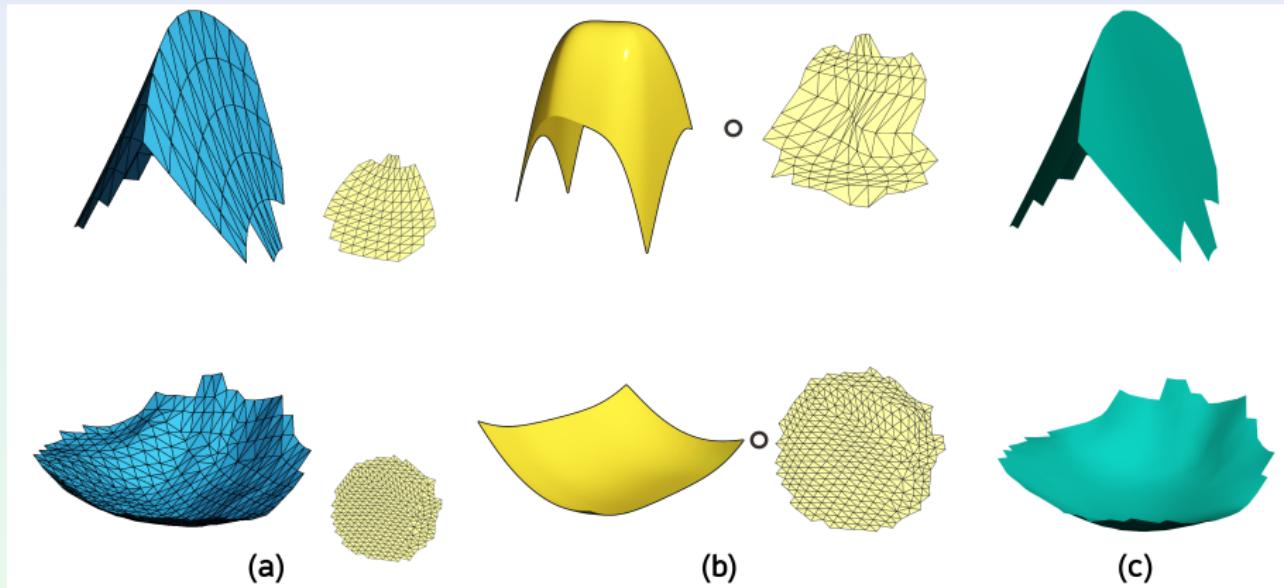
$$\begin{aligned} \min_{c, U} \quad & F(c, U) + \alpha E(c, U) + Q(U) \\ \text{s.t.} \quad & \|c\|_0 \leq s \end{aligned} \tag{37}$$

where  $F(c, U)$  is the data term,  $E(c, U)$  is a geometric smoothness term decided by both coordinate  $c$  and  $U$ , and  $Q(U)$  is the regularization term of the triangle structure of  $\{u_i\}$ .

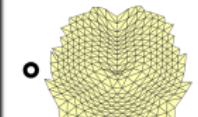
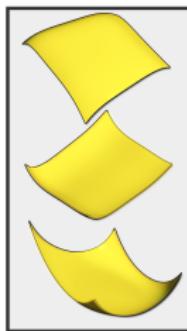
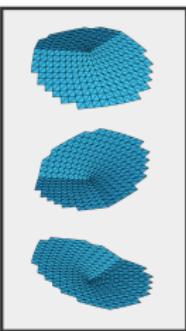
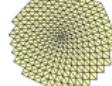
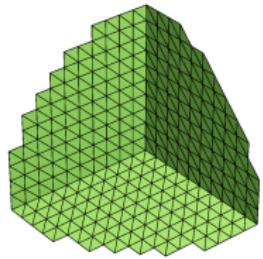
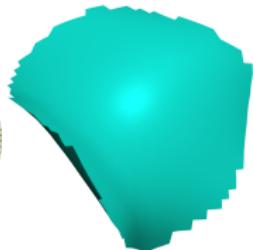
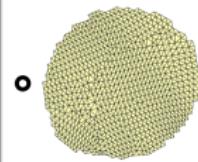
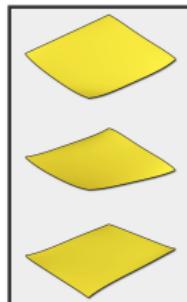
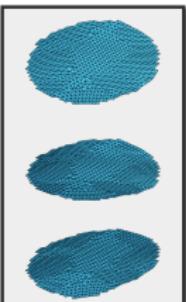
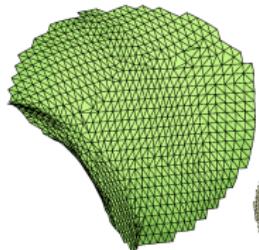
# Overview of our algorithm



## Experimental results



# Experimental results



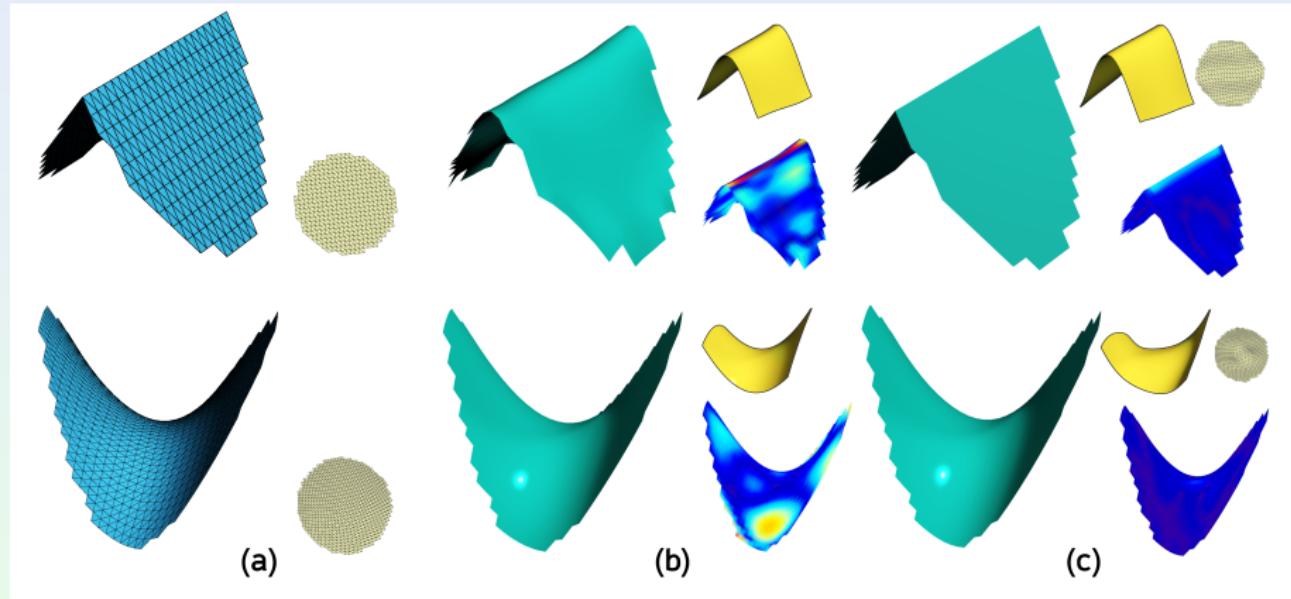
(a)

(b)

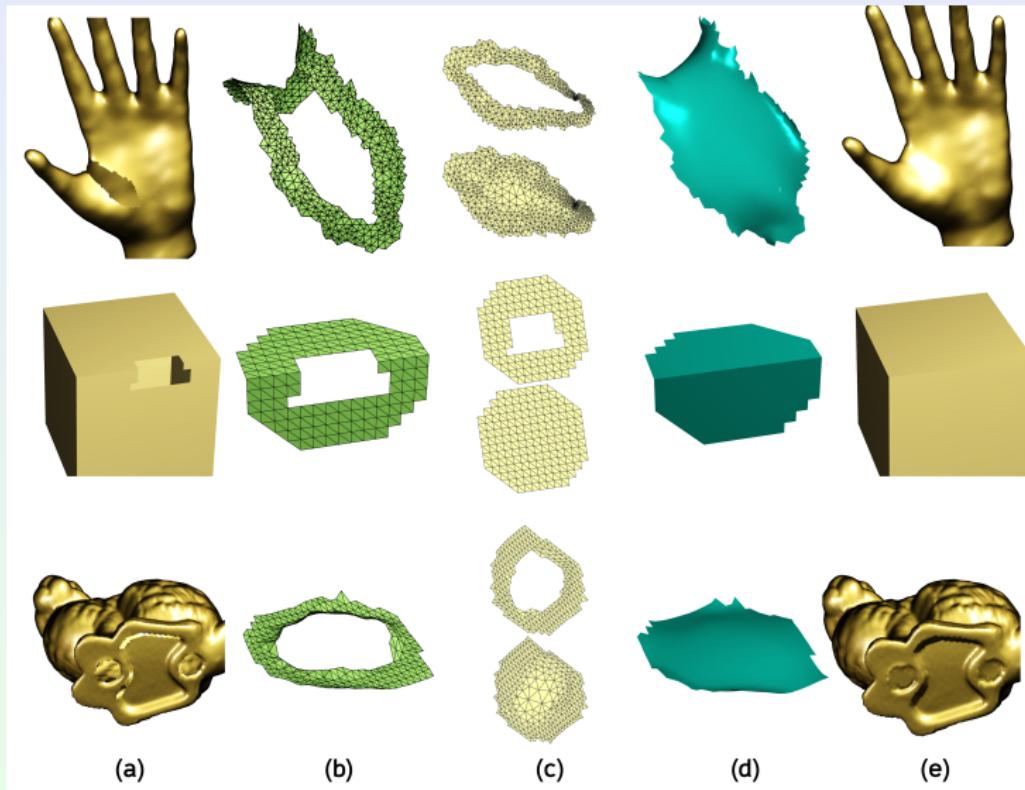
(c)

(d)

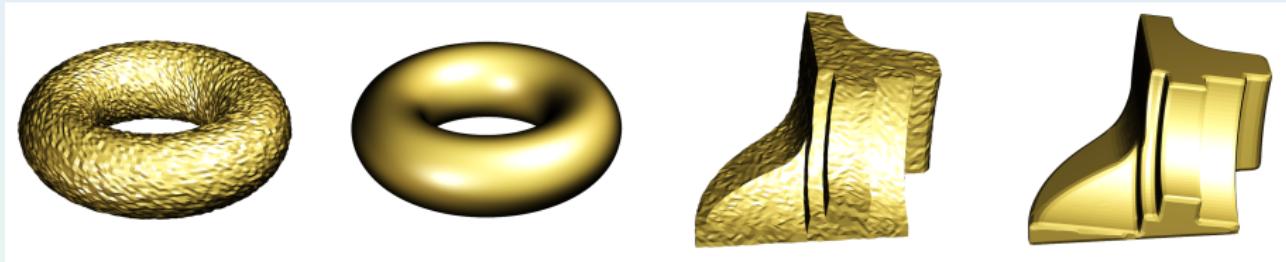
## Comparison with naive sparse fitting



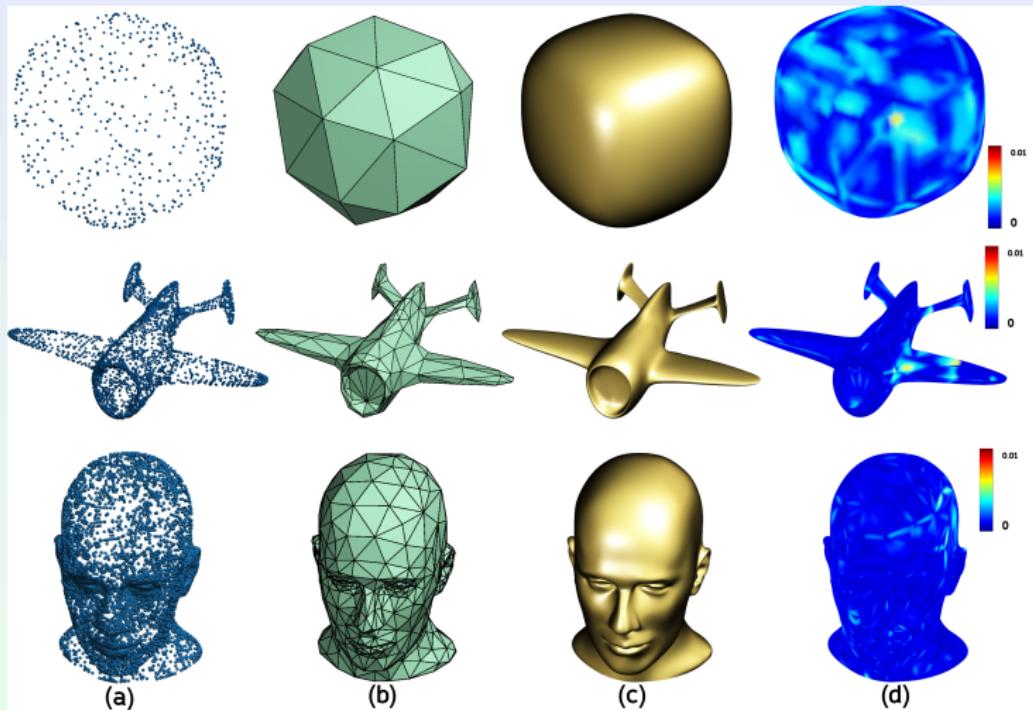
# Application to mesh inpainting



# Application to mesh smoothing



# Application to surface reconstruction of point cloud



# Outline I

- ① Compressed Sensing
- ② Sparse Modeling
  - Sparsity-seeking representations
  - Numerical optimization
  - Applications
- ③ Sparse Optimization
  - Sparse Optimization Models
  - Sparse Optimization Algorithms
- ④ Decoupling Noises and Features via  $\ell_1$ -analysis Compressed Sensing
  - Discrete Laplacian regularization smoothing
  - Feature recovering via  $\ell_1$ -analysis optimization

# Outline II

- Experimental results

5 Construction of Manifolds via Consistent Sparse Representations

6 Sparse Representation with Parameterization Optimization

7 Conclusion and Future Work

## Future plans . . .

We would like to explore the emerging framework of sparse modeling and sparse optimization to interdisciplinary applications.

Thanks for your attention!