

Compressed Sensing Meets Machine Learning

- Classification of Mixture Subspace Models via Sparse Representation

Allen Y. Yang
<yang@eecs.berkeley.edu>

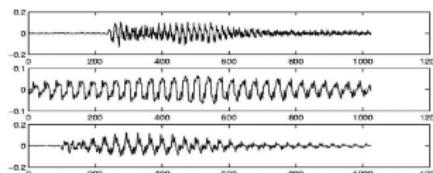
Feb. 25, 2008. UC Berkeley



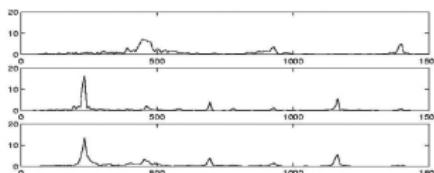
What is Sparsity

Sparsity

A signal is sparse if most of its coefficients are (approximately) zero.



(a) Harmonic functions



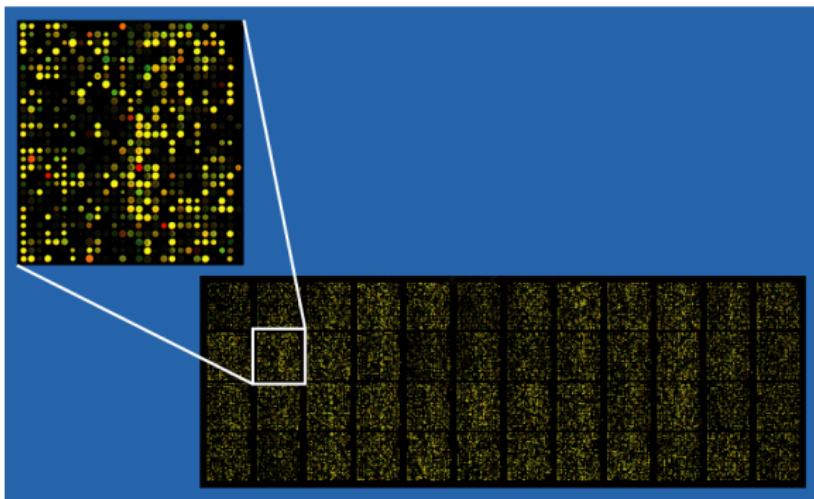
(b) Magnitude spectrum



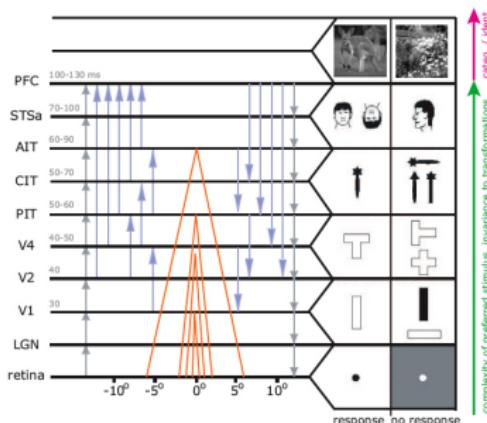
Figure: 2-D DCT transform.

Sparsity in spatial domain

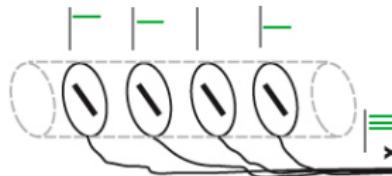
- gene microarray data [Drmanac et al. 1993]



- Sparsity in human visual cortex [Olshausen & Field 1997, Serre & Poggio 2006]



- ① **Feed-forward:** No iterative feedback loop.
- ② **Redundancy:** Average 80-200 neurons for each feature representation.
- ③ **Recognition:** Information exchange between stages is not about individual neurons, but rather **how many neurons as a group fire together**.



Sparsity and ℓ^1 -Minimization

- “Black gold” age [Claerbout & Muir 1973, Taylor, Banks & McCoy 1979]

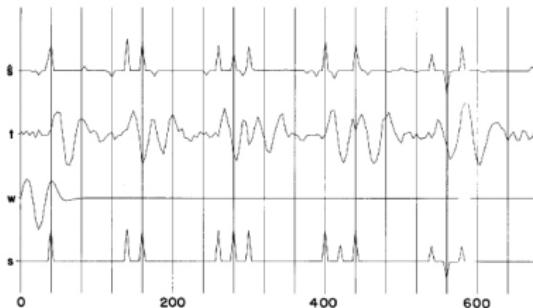
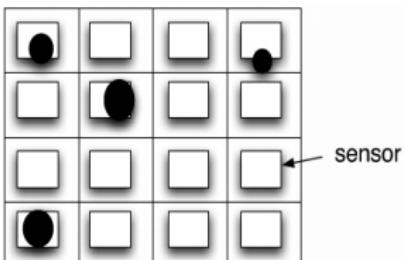


Figure: Deconvolution of spike train.

Sparse Support Estimators

- **Sparse support estimator** [Donoho 1992, Meinshausen & Bühlmann 2006, Yu 2006, Wainwright 2006, Ramchandran 2007, Gastpar 2007]



- **Basis pursuit** [Chen & Donoho 1999]: Given $\mathbf{y} = A\mathbf{x}$ and \mathbf{x} unknown,

$$\mathbf{x}^* = \arg \min \|\mathbf{x}\|_1, \text{ subject to } \mathbf{y} = A\mathbf{x}$$

- **The Lasso** (least absolute shrinkage and selection operator) [Tibshirani 1996]

$$\mathbf{x}^* = \arg \min \|\mathbf{y} - A\mathbf{x}\|_2, \text{ subject to } \|\mathbf{x}\|_1 \leq k$$

Taking Advantage of Sparsity

What generates sparsity? (*d'après Emmanuel Candès*)

- Measure first, analyze later.
- Curse of dimensionality.

- ① Numerical analysis: sparsity reduces cost for storage and computation.
- ② Regularization in classification:

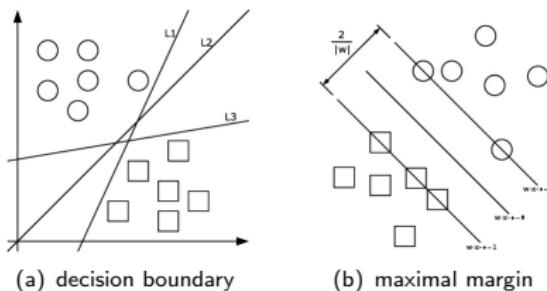


Figure: Linear support vector machine (SVM)

Our Contributions

- ① Classification via compressed sensing
- ② Performance in face recognition
- ③ Extensions
 - Outlier rejection
 - Occlusion compensation
- ④ Distributed pattern recognition in sensor networks.

Problem Formulation in Face Recognition

① Notations

- Training: For K classes, collect training samples $\{\mathbf{v}_{1,1}, \dots, \mathbf{v}_{1,n_1}\}, \dots, \{\mathbf{v}_{K,1}, \dots, \mathbf{v}_{K,n_K}\} \in \mathbb{R}^D$.
 - Test: Present a new $\mathbf{y} \in \mathbb{R}^D$, solve for label(\mathbf{y}) $\in [1, 2, \dots, K]$.

② Construct \mathbb{R}^D sample space via **stacking**

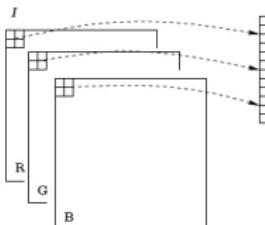
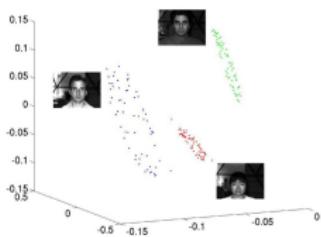


Figure: For images, assume 3-channel 640×480 image, $D = 3 \cdot 640 \cdot 480 \approx 1\text{e}6$.

③ Assume y belongs to Class i [Belhumeur et al. 1997, Basri & Jacobs 2003]



$$\mathbf{y} = \alpha_{i,1}\mathbf{v}_{i,1} + \alpha_{i,2}\mathbf{v}_{i,2} + \cdots + \alpha_{i,n_i}\mathbf{v}_{i,n_i},$$

$$= A_i\alpha_i,$$

where $A_i = [\mathbf{v}_{i,1}, \mathbf{v}_{i,2}, \dots, \mathbf{v}_{i,n_i}]$:

① Nevertheless, i is the variable we need to solve.

Global representation:

$$\begin{aligned}\mathbf{y} &= [A_1, A_2, \dots, A_K] \begin{bmatrix} \alpha_1 \\ \alpha_2 \\ \vdots \\ \alpha_K \end{bmatrix}, \\ &= A\mathbf{x}_0.\end{aligned}$$

② Over-determined system: $A \in \mathbb{R}^{D \times n}$, where $D \gg n = n_1 + \dots + n_K$.

\mathbf{x}_0 encodes membership of \mathbf{y} : If \mathbf{y} belongs to Subject i ,

$$\mathbf{x}_0 = [0 \dots 0 \ \alpha_i \ 0 \dots 0]^T \in \mathbb{R}^n.$$

Problems to face

- Solving for \mathbf{x}_0 in \mathbb{R}^D is **intractable**.
- True solution \mathbf{x}_0 is **sparse**: Average $\frac{1}{K}$ terms non-zero.

Dimensionality Reduction

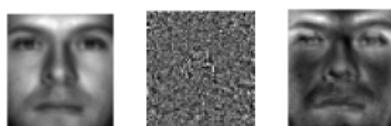
- ① Construct linear projection $R \in \mathbb{R}^{d \times D}$, d is the **feature dimension**, $d \ll D$.

$$\tilde{\mathbf{y}} \doteq R\mathbf{y} = RA\mathbf{x}_0 = \tilde{A}\mathbf{x}_0 \in \mathbb{R}^d.$$

$\tilde{A} \in \mathbb{R}^{d \times n}$, but \mathbf{x}_0 is unchanged.

- ② Holistic features

- Eigenfaces [Turk 1991]
- Fisherfaces [Belhumeur 1997]
- Laplacianfaces [He 2005]

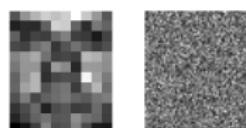


- ③ Partial features



- ④ Unconventional features

- Downsampled faces
- Random projections



ℓ^0 -Minimization

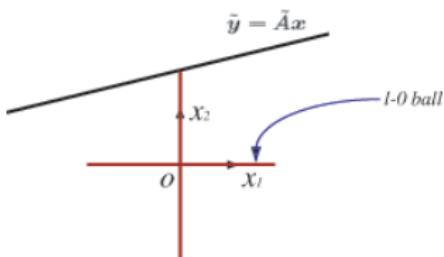
① Solving for **sparsest** solution via ℓ^0 -Minimization

$$\mathbf{x}_0 = \arg \min_{\mathbf{x}} \|\mathbf{x}\|_0 \text{ s.t. } \tilde{\mathbf{y}} = \tilde{\mathbf{A}}\mathbf{x}.$$

$\|\cdot\|_0$ simply counts the number of nonzero terms.

② ℓ^0 -Ball

- ℓ^0 -ball is not convex.
- ℓ^0 -minimization is NP-hard.



ℓ^1/ℓ^0 Equivalence

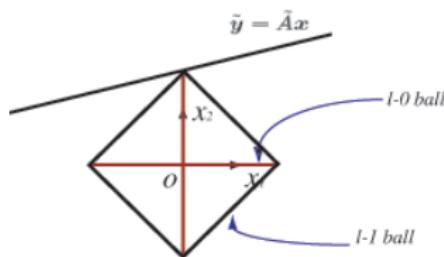
① **Compressed sensing:** If x_0 is sparse enough, ℓ^0 -minimization is equivalent to

$$(P_1) \quad \min \|x\|_1 \text{ s.t. } \tilde{y} = \tilde{A}x.$$

$$\|x\|_1 = |x_1| + |x_2| + \cdots + |x_n|.$$

② ℓ^1 -Ball

- ℓ^1 -Minimization is convex.
- Solution equal to ℓ^0 -minimization.



③ ℓ^1/ℓ^0 Equivalence: [Donoho 2002, 2004; Candes et al. 2004; Baraniuk 2006]

Given $\tilde{y} = \tilde{A}x_0$, there exists **equivalence breakdown point** (EBP) $\rho(\tilde{A})$, if $\|x_0\|_0 < \rho$:

- ℓ^1 -solution is unique
- $x_1 = x_0$

ℓ^1 -Minimization Routines

• Matching pursuit [Mallat 1993]

- ① Find most correlated vector \mathbf{v}_i in $\tilde{\mathbf{A}}$ with \mathbf{y} : $i = \arg \max \langle \mathbf{y}, \mathbf{v}_j \rangle$.
- ② $\tilde{\mathbf{A}} \leftarrow \tilde{\mathbf{A}}^i$, $\mathbf{x}_i \leftarrow \langle \mathbf{y}, \mathbf{v}_i \rangle$, $\mathbf{y} \leftarrow \mathbf{y} - \mathbf{x}_i \mathbf{v}_i$.
- ③ Repeat until $\|\mathbf{y}\| < \epsilon$.

• Basis pursuit [Chen 1998]

- ① Assume \mathbf{x}_0 is m -sparse.
- ② Select m linearly independent vectors B_m in $\tilde{\mathbf{A}}$ as a basis

$$\mathbf{x}_m = B_m^\dagger \mathbf{y}.$$

- ③ Repeat swapping one basis vector in B_m with another vector in $\tilde{\mathbf{A}}$ if improve $\|\mathbf{y} - B_m \mathbf{x}_m\|$.
- ④ If $\|\mathbf{y} - B_m \mathbf{x}_m\|_2 < \epsilon$, stop.

• Quadratic solvers: $\tilde{\mathbf{y}} = \tilde{\mathbf{A}} \mathbf{x}_0 + \mathbf{z} \in \mathbb{R}^d$, where $\|\mathbf{z}\|_2 < \epsilon$

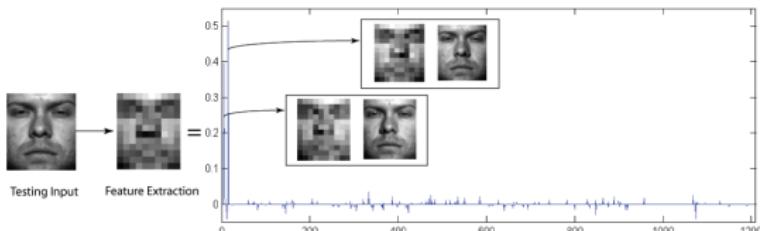
$$\mathbf{x}^* = \arg \min \{ \|\mathbf{x}\|_1 + \lambda \|\mathbf{y} - \tilde{\mathbf{A}} \mathbf{x}\|_2 \}$$

[Lasso, Second-order cone programming]: More expensive.

Matlab Toolboxes

- **ℓ^1 -Magic** by Candès at Caltech.
- **SparseLab** by Donoho at Stanford.
- **cvx** by Boyd at Stanford.

Classification

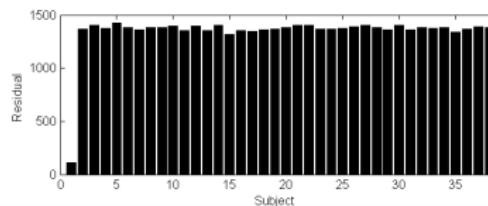


① Project \mathbf{x}_1 onto face subspaces:

$$\delta_1(\mathbf{x}_1) = \begin{bmatrix} \alpha_1 \\ 0 \\ \vdots \\ 0 \end{bmatrix}, \delta_2(\mathbf{x}_1) = \begin{bmatrix} 0 \\ \alpha_2 \\ \vdots \\ 0 \end{bmatrix}, \dots, \delta_K(\mathbf{x}_1) = \begin{bmatrix} 0 \\ 0 \\ \vdots \\ \alpha_K \end{bmatrix}. \quad (1)$$

② Define residual $r_i = \|\tilde{\mathbf{y}} - \tilde{\mathbf{A}}\delta_i(\mathbf{x}_1)\|_2$ for Subject i :

- $\text{id}(\mathbf{y}) = \arg \min_{i=1, \dots, K} \{r_i\}$



AR Database 100 Subjects (Illumination and Expression Variance)

**Table: I.** Nearest Neighbor

Dimension	30	54	130	540
Eigen [%]	68.1	74.8	79.3	80.5
Laplacian [%]	73.1	77.1	83.8	89.7
Random [%]	56.7	63.7	71.4	75
Down [%]	51.7	60.9	69.2	73.7
Fisher [%]	83.4	86.8	N/A	N/A

Table: II. Nearest Subspace

	30	54	130	540
64.1	77.1	82	85.1	
66	77.5	84.3	90.3	
59.2	68.2	80	83.3	
56.2	67.7	77	82.1	
80.3	85.8	N/A	N/A	

Table: III. Linear SVM

Dimension	30	54	130	540
Eigen [%]	73	84.3	89	92
Laplacian [%]	73.4	85.8	90.8	95.7
Random [%]	54.1	70.8	81.6	88.8
Down [%]	51.4	73	83.4	90.3
Fisher [%]	86.3	93.3	N/A	N/A

Table: IV. ℓ^1 -Minimization

	30	54	130	540
71.1	80	85.7	92	
73.7	84.7	91	94.3	
57.8	75.5	87.6	94.7	
46.8	67	84.6	93.9	
87	92.3	N/A	N/A	

Sparsity vs. Non-sparsity: ℓ^1 and SVM decisively outperform NN and NS.

- ① Our framework seeks sparsity in representation of \mathbf{y} .
- ② SVM seeks sparsity in decision boundaries on $A = [\mathbf{v}_1, \dots, \mathbf{v}_n]$.
- ③ NN and NS do not enforce sparsity.

ℓ^1 -Minimization vs. SVM: Performance of SVM depends on the choice of features.

- ① Random project performs poorly with SVMs.
- ② ℓ^1 -Minimization guarantees performance convergence with different features.
- ③ At lower-dimensional space, *Fisher* features outperform.

Table: III. Linear SVM

Dimension	30	54	130	540
Eigen [%]	73	84.3	89	92
Laplacian [%]	73.4	85.8	90.8	95.7
Random [%]	54.1	70.8	81.6	88.8
Down [%]	51.4	73	83.4	90.3
Fisher [%]	86.3	93.3	N/A	N/A

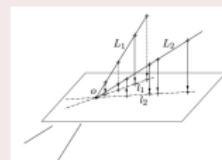
Table: IV. ℓ^1 -Minimization

30	54	130	540
71.1	80	85.7	92
73.7	84.7	91	94.3
57.8	75.5	87.6	94.7
46.8	67	84.6	93.9
87	92.3	N/A	N/A

Randomfaces

Blessing of Dimensionality [Donoho 2000]

In high-dimensional data space \mathbb{R}^D , with **overwhelming probability**,
 ℓ^1/ℓ^0 equivalence holds for random projection R .



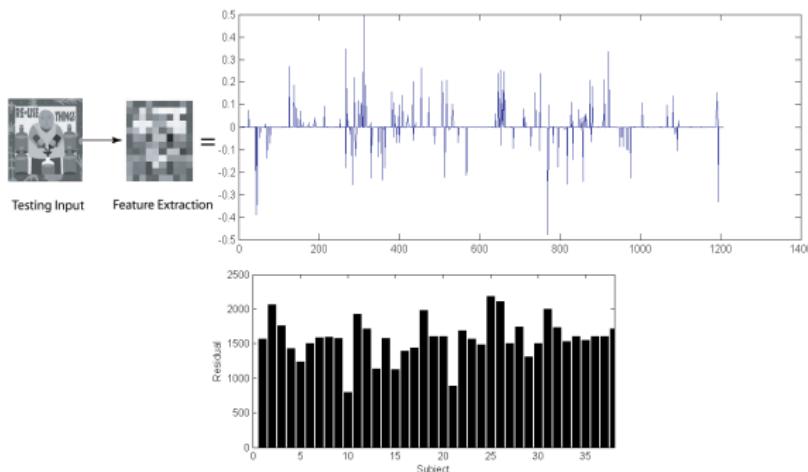
Unconventional properties:

- ① Domain independent!
- ② Data independent!
- ③ Fast to generate and compute!

Reference: Yang et al. *Feature selection in face recognition: A sparse representation perspective*. Berkeley Tech Report, 2007.

Variation: Outlier Rejection

- ℓ^1 -Coefficients for invalid images

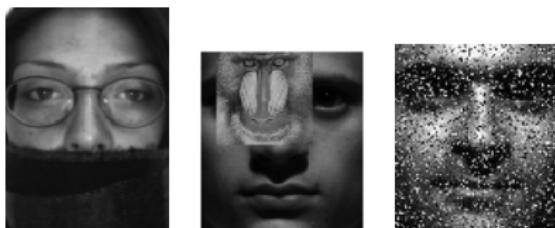


Outlier Rejection

When ℓ^1 -solution is not sparse or concentrated to one subspace, the test sample is invalid.

$$\text{Sparsity Concentration Index: } \text{SCI}(\mathbf{x}) \doteq \frac{K \cdot \max_i \|\delta_i(\mathbf{x})\|_1 / \|\mathbf{x}\|_1 - 1}{K - 1} \in [0, 1].$$

Variation: Occlusion Compensation



- ① Sparse representation + sparse error

$$\mathbf{y} = \mathbf{Ax} + \mathbf{e}$$



- ② Occlusion compensation:

$$\mathbf{y} = (\mathbf{A} \quad | \quad \mathbf{I}) \begin{pmatrix} \mathbf{x} \\ \mathbf{e} \end{pmatrix} = \mathbf{Bw}$$

Reference: Wright et al. *Robust face recognition via sparse representation*. UIUC Tech Report, 2007.

Distributed Pattern Recognition

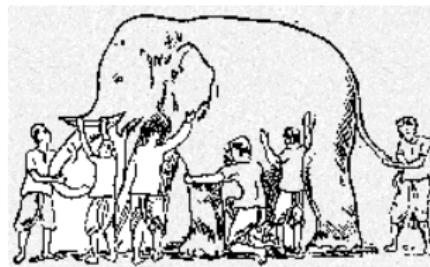


Figure: d-Oracle: Distributed object recognition via camera wireless networks.

Key components:

- ① Each sensor only observes partial profile of the event: Demand a global classification framework.
- ② Individual sensor obtains limited classification ability: Sensors become active only when certain events are locally detected.
- ③ The network configuration is dynamic: Global classifier needs to adapt to change of active sensors.

Problem Formulation for Distributed Action Recognition

Architecture

- 8 sensors distributed on human body.
- Location of the sensors are given and fixed.
- Each sensor carries triaxial accelerometer and biaxial gyroscope.
- Sampling frequency: 20Hz.

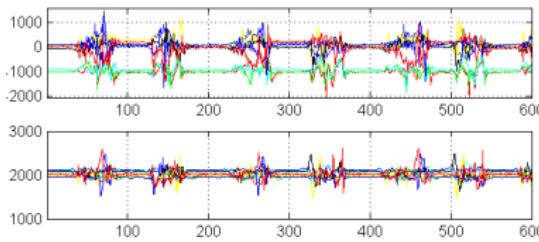
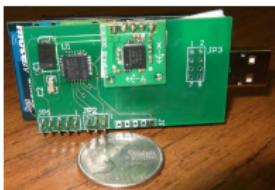


Figure: Readings from 8 x-axis accelerometers and x-axis gyroscopes for a *stand-kneel-stand* sequence.

Challenges for Distributed Action Recognition

- ① Simultaneous segmentation and classification.
- ② Individual sensors **not sufficient** to classify all human actions.
- ③ Simulate **sensor failure and network congestion** by different subsets of active sensors.
- ④ **Identity independence:** The prior examples of the subject for testing are excluded as part of training data.

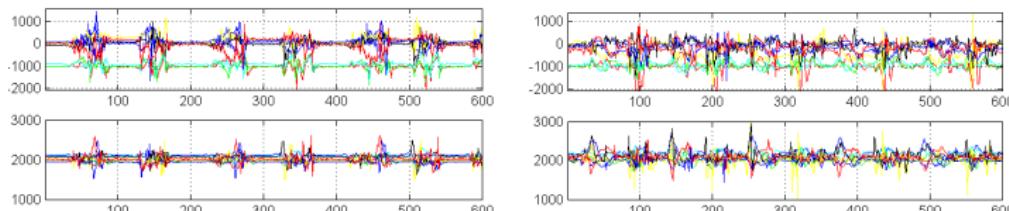
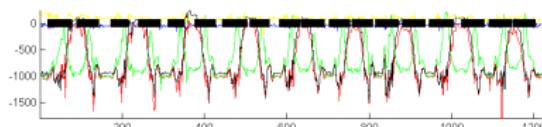


Figure: Same actions performed by two subjects.

Mixture Subspace Model for Distributed Action Recognition

- ① Training samples are segmented manually with correct labels.



- ② On each sensor node i , normalize the vector form (via stacking)

$$\mathbf{v}_i = [x(1), \dots, x(h), y(1), \dots, y(h), z(1), \dots, z(h), \theta(1), \dots, \theta(h), \rho(1), \dots, \rho(h)]^T \in \mathbb{R}^{5h}$$

- ③ Full body motion

$$\text{Training sample: } \mathbf{v} = \begin{pmatrix} \mathbf{v}_1 \\ \vdots \\ \mathbf{v}_8 \end{pmatrix} \quad \text{Test sample: } \mathbf{y} = \begin{pmatrix} \mathbf{y}_1 \\ \vdots \\ \mathbf{y}_8 \end{pmatrix} \in \mathbb{R}^{8 \cdot 5h}$$

- ④ Mixture subspace model

$$\mathbf{y} = \begin{pmatrix} \mathbf{y}_1 \\ \vdots \\ \mathbf{y}_8 \end{pmatrix} = \left(\begin{pmatrix} \mathbf{v}_1 \\ \vdots \\ \mathbf{v}_8 \end{pmatrix}_1, \dots, \begin{pmatrix} \mathbf{v}_1 \\ \vdots \\ \mathbf{v}_8 \end{pmatrix}_n \right) \mathbf{x} = \mathbf{Ax}.$$

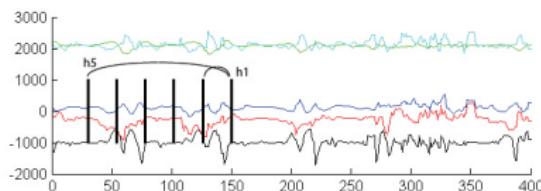
Localized Classifiers

Distributed Sparse Representation

$$\begin{pmatrix} \mathbf{y}_1 \\ \vdots \\ \mathbf{y}_8 \end{pmatrix} = \left(\begin{pmatrix} \mathbf{v}_1 \\ \vdots \\ \mathbf{v}_8 \end{pmatrix}_1, \dots, \begin{pmatrix} \mathbf{v}_1 \\ \vdots \\ \mathbf{v}_8 \end{pmatrix}_n \right) \mathbf{x} \Leftrightarrow \begin{cases} \mathbf{y}_1 = (\mathbf{v}_{1,1}, \dots, \mathbf{v}_{1,n}) \mathbf{x} \\ \vdots \\ \mathbf{y}_8 = (\mathbf{v}_{8,1}, \dots, \mathbf{v}_{8,n}) \mathbf{x} \end{cases}$$

On each sensor node i :

- ① Given a (long) test sequence at time t , apply multiple duration hypotheses: $\mathbf{y}_i \in \mathbb{R}^{5h}$.



- ② Choose Fisher features $R_i \in \mathbb{R}^{10 \times 5h}$:

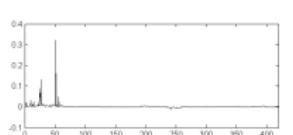
$$\tilde{\mathbf{y}}_i = R_i \mathbf{y}_i = R_i A_i \mathbf{x} = \tilde{A}_i \mathbf{x} \in \mathbb{R}^{10}$$

Localized Classifiers

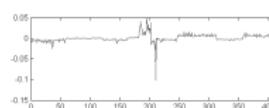
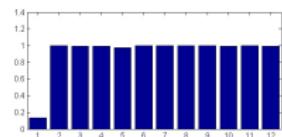
- ① Equivalently, **define** $R'_i = (0 \cdots R_i \cdots 0)$

$$\tilde{\mathbf{y}}_i = (0 \cdots R_i \cdots 0) \begin{pmatrix} y_1 \\ \vdots \\ y_8 \end{pmatrix} = (0 \cdots R_i \cdots 0) \left(\left(\begin{pmatrix} \mathbf{v}_1 \\ \vdots \\ \mathbf{v}_8 \end{pmatrix} \right)_1, \dots, \left(\begin{pmatrix} \mathbf{v}_1 \\ \vdots \\ \mathbf{v}_8 \end{pmatrix} \right)_n \right) \mathbf{x} = R'_i A \mathbf{x} \in \mathbb{R}^{10}$$

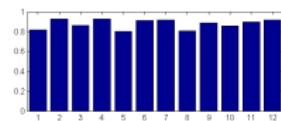
- ② For all segmentation hypotheses, apply **sparsity concentration index (SCI)** threshold σ_1 :



(a) valid segmentation



(b) invalid segmentation



Local Sparsity Threshold σ_1

- If $\text{SCI}(\mathbf{x}) > \sigma_1$, sensor i becomes **active** and transmits $\tilde{\mathbf{y}}_i \in \mathbb{R}^{10}$.
- $\tilde{\mathbf{y}}_i$ provides a **segmentation hypothesis** at time t and length h_i .

Adaptive Global Classifier

- Adaptive classification for a subset of active sensors (Suppose $1, \dots, L$ at time t and h_i)

Define **global feature matrix** $R' = \begin{pmatrix} R_1 & \cdots & 0 & \cdots & 0 \\ \vdots & \ddots & \vdots & & \vdots \\ 0 & \cdots & R_L & \cdots & 0 \end{pmatrix}$:

$$\begin{pmatrix} \tilde{y}_1 \\ \vdots \\ \tilde{y}_L \end{pmatrix} = R' \begin{pmatrix} y_1 \\ \vdots \\ y_8 \end{pmatrix} = R' \begin{pmatrix} A_1 \\ \vdots \\ A_8 \end{pmatrix} \mathbf{x} = R' \mathbf{A} \mathbf{x}$$

- Global segmentation: Regardless of L active sensors, given a global threshold σ_2 ,
If $\text{SCI}(\mathbf{x}) > \sigma_2$, accept \mathbf{y} as **global segmentation with label given by \mathbf{x}** .

Distributed Classification via Compressed Sensing

- Reformulate adaptive classification via feature matrix R' :

Local: $R' = (0 \cdots R_i \cdots 0)$ Global: $R' = \begin{pmatrix} R_1 & \cdots & 0 & \cdots & 0 \\ \vdots & \ddots & \vdots & & \vdots \\ 0 & \cdots & R_L & \cdots & 0 \end{pmatrix} \Leftrightarrow R' \mathbf{y} = R' \mathbf{A} \mathbf{x}$

- The representation \mathbf{x} and training matrix A remain **invariant**.
- Segmentation, recognition, and outlier rejection are unified on \mathbf{x}** .

Experiment

- Algorithm only depends on two parameters: (σ_1, σ_2) .
- Data communications between sensors and station are 10-D action features.

Precision vs Recall:

Sensors	2	7	2,7	1,2,7	1- 3, 7,8	1- 8
Prec [%]	89.8	94.6	94.4	92.8	94.6	98.8
Rec [%]	65	61.5	82.5	80.6	89.5	94.2



Confusion Tables:

Class (total)	1	2	3	4	5	6	7	8	9	10	11	12
1 StSi (60)	60	0	0	0	0	0	0	0	0	0	0	0
2 SiSt (60)	0	52	0	0	0	0	0	0	0	0	0	0
3 SiLi (62)	1	0	58	0	0	0	0	0	0	0	0	0
4 LiSi (62)	0	0	0	60	0	0	0	0	0	0	0	0
5 Bend (30)	1	0	0	0	29	0	0	0	0	0	0	0
6 StKn (33)	0	0	0	0	31	0	0	0	0	0	0	0
7 KnSt (30)	0	0	0	0	0	30	0	0	0	1	0	0
8 RoR (95)	0	0	0	0	0	0	93	0	0	0	1	0
9 RoL (96)	0	0	0	0	0	0	0	96	0	0	0	0
10 Jump (34)	0	0	0	0	0	0	0	0	31	0	0	0
11 Up (33)	0	0	0	0	0	0	0	0	0	24	0	0
12 Down (31)	0	0	0	0	0	0	0	0	0	0	3	26

(a) Sensor 1-8

Class (total)	1	2	3	4	5	6	7	8	9	10	11	12
1 StSi (60)	0	0	0	2	0	0	0	0	0	0	0	0
2 SiSt (60)	0	1	0	0	1	0	0	2	0	0	0	0
3 SiLi (62)	0	0	59	0	0	0	0	0	0	0	0	0
4 LiSi (62)	0	0	0	60	0	0	0	0	0	0	0	0
5 Bend (30)	4	0	0	0	9	0	0	0	0	0	0	0
6 StKn (33)	0	0	0	0	0	28	0	0	0	0	0	0
7 KnSt (30)	0	0	0	0	0	0	20	0	1	0	0	0
8 RoR (95)	0	0	0	0	1	0	0	82	0	0	0	0
9 RoL (96)	0	0	0	0	0	0	0	0	90	0	0	0
10 Jump (34)	0	0	0	0	0	0	0	0	0	22	0	0
11 Up (33)	1	0	0	0	3	0	0	1	0	0	4	0
12 Down (31)	0	0	0	0	0	0	0	0	5	0	0	10

(b) Sensor 7

Reference: Yang et al. *Distributed Segmentation and Classification of Human Actions Using a Wearable Motion Sensor Network*.
Berkeley Tech Report, 2007.

Conclusion

- ① Sparsity is important for classification of HD data.
- ② A new recognition framework via **compressed sensing**.
- ③ In HD feature space, choosing an “optimal” feature becomes not significant.
- ④ **Randomfaces, outliers, occlusion.**
- ⑤ Distributed pattern recognition in body sensor networks.

Future Directions

- Distributed camera networks



- Biosensor networks in health care



Acknowledgments

Collaborators

- Berkeley: Shankar Sastry, Ruzena Bajcsy
- UIUC: Yi Ma
- UT-Dallas: Roozbeh Jafari

MATLAB Toolboxes

- **ℓ^1 -Magic** by Candès at Caltech.
- **SparseLab** by Donoho at Stanford.
- **cvx** by Boyd at Stanford.

References

- *Robust face recognition via sparse representation*. Submitted to PAMI, 2008.
- *Distributed segmentation and classification of human actions using a wearable motion sensor network*. Berkeley Tech Report 2007.