

# Discriminative Graphical Models for Sparsity-Based Hyperspectral Target Detection

Umamahesh Srinivas<sup>1</sup>

Yi Chen<sup>2</sup>

Vishal Monga<sup>1</sup>

Nasser Nasrabadi<sup>3</sup>

Trac Tran<sup>2</sup>

<sup>1</sup>Pennsylvania State University, University Park, USA

<sup>2</sup>The Johns Hopkins University, Baltimore, USA

<sup>3</sup>U.S. Army Research Laboratory, Adelphi, USA



**IEEE International Geoscience and Remote Sensing Symposium**

July 24, 2012

# Outline

- 1 Overview: Hyperspectral target detection
- 2 Sparse representation-based target detection
- 3 Contribution: Discriminative graphical models for target detection
- 4 Experiments and results

# Hyperspectral Imaging (HSI)

- Materials reflect, absorb, and emit electromagnetic energy at different wavelengths in a specific manner
- **Imaging spectrometer** measures reflectance over many contiguous spectral wavelength bands
- **Applications:** military, agriculture, mineralogy, etc.

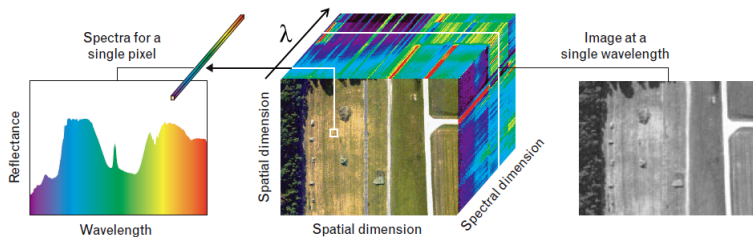


Figure : Illustration of a HSI data cube<sup>1</sup>.

- Spectral range: 400 nm - 2500 nm (visible to near infra-red)

<sup>1</sup>Manolakis et al., Lincoln Lab. Journal, 2003

# Hyperspectral target detection

- Binary hypothesis testing problem (at each pixel independently):

$H_0$  : target absent

$H_1$  : target present.

- Likelihood ratio test:

$$L(\mathbf{x}) = \frac{p(\mathbf{x}|H_1)}{p(\mathbf{x}|H_0)} \underset{H_0}{\overset{H_1}{\gtrless}} \tau.$$

- Issues:

- Spectral variability  $\rightarrow$  atmospheric conditions, sensor noise, material composition
- Disproportionately small number of target pixels in scene

# HSI Target detection: Prior work

- Adaptive matched filter (AMF)<sup>2</sup>

$$H_0 : \mathbf{x} = \mathbf{n}$$

$$H_1 : \mathbf{x} = a\mathbf{s} + \mathbf{n}.$$

$$L(\mathbf{x}) = \frac{\mathbf{s}^T \hat{\mathbf{C}}_n^{-1} \mathbf{x}}{\mathbf{s}^T \hat{\mathbf{C}}_n^{-1} \mathbf{s}},$$

where  $\mathbf{s} \rightarrow$  target spectral signature,  $\hat{\mathbf{C}}_n \rightarrow$  background covariance.

---

<sup>2</sup>Robey et al., IEEE Trans. Aerosp. Electron. Syst., 1992

<sup>3</sup>Scharf et al., IEEE Trans. Signal Process., 1994

# HSI Target detection: Prior work

- Adaptive matched filter (AMF)<sup>2</sup>

$$H_0 : \mathbf{x} = \mathbf{n}$$

$$H_1 : \mathbf{x} = a\mathbf{s} + \mathbf{n}.$$

$$L(\mathbf{x}) = \frac{\mathbf{s}^T \hat{\mathbf{C}}_n^{-1} \mathbf{x}}{\mathbf{s}^T \hat{\mathbf{C}}_n^{-1} \mathbf{s}},$$

where  $\mathbf{s} \rightarrow$  target spectral signature,  $\hat{\mathbf{C}}_n \rightarrow$  background covariance.

- Matched subspace detector (MSD)<sup>3</sup>

$$H_0 : \mathbf{x} = \mathbf{B}\boldsymbol{\zeta} + \mathbf{n}$$

$$H_1 : \mathbf{x} = \mathbf{T}\boldsymbol{\theta} + \mathbf{B}\boldsymbol{\zeta} + \mathbf{n}.$$

$$L(\mathbf{x}) = \frac{\mathbf{x}^T \mathbf{P}_B^\perp \mathbf{x}}{\mathbf{x}^T \mathbf{P}_{TB}^\perp \mathbf{x}},$$

where  $\mathbf{P}_B^\perp = \mathbf{I} - \mathbf{P}_B$  and  $\mathbf{P}_{TB}^\perp = \mathbf{I} - \mathbf{P}_{TB}$  are projection matrices.

---

<sup>2</sup>Robey et al., IEEE Trans. Aerosp. Electron. Syst., 1992

<sup>3</sup>Scharf et al., IEEE Trans. Signal Process., 1994

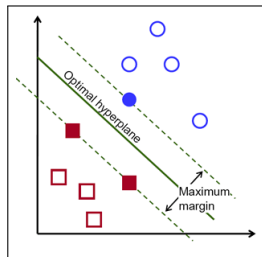
# HSI Target detection and classification: Prior work

- Support vector machines (SVM)<sup>4</sup>

Margin-maximizing hyperplane:

$$f(\mathbf{x}) = \sum_{i=1}^n \alpha_i y_i K(\mathbf{s}_i, \mathbf{x}) + \beta,$$

where  $\mathbf{s}_i \rightarrow$  support vectors,  
 $y_i \in \{-1, +1\}$ ,  $K \rightarrow$  kernel.



<sup>4</sup>Melgani et al., IEEE Trans. Geosci. Remote Sens., 2004

<sup>5</sup>Camps-Valls et al., IEEE Geosci. Remote Sens. Lett., 2006

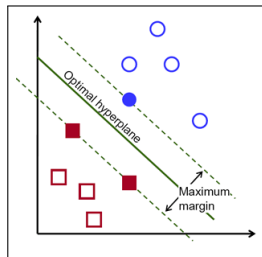
# HSI Target detection and classification: Prior work

- Support vector machines (SVM)<sup>4</sup>

Margin-maximizing hyperplane:

$$f(\mathbf{x}) = \sum_{i=1}^n \alpha_i y_i K(\mathbf{s}_i, \mathbf{x}) + \beta,$$

where  $\mathbf{s}_i \rightarrow$  support vectors,  
 $y_i \in \{-1, +1\}$ ,  $K \rightarrow$  kernel.



- Composite kernel SVM<sup>5</sup>
  - Fusion of spatial and spectral information

$$K(\mathbf{x}_i, \mathbf{x}_j) = \mu K_s(\mathbf{x}_i^s, \mathbf{x}_j^s) + (1 - \mu) K_w(\mathbf{x}_i^w, \mathbf{x}_j^w),$$

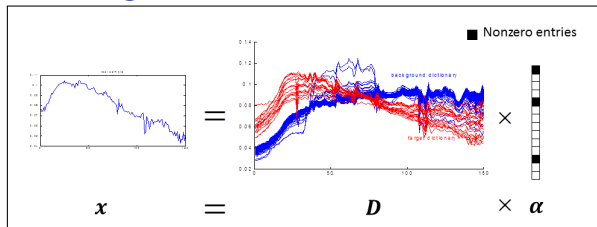
where  $\mu \in [0, 1]$ ,  $x_i^w \rightarrow$  spectral pixel,  $x_i^s \rightarrow$  spatial feature from a local neighborhood.

<sup>4</sup>Melgani et al., IEEE Trans. Geosci. Remote Sens., 2004

<sup>5</sup>Camps-Valls et al., IEEE Geosci. Remote Sens. Lett., 2006

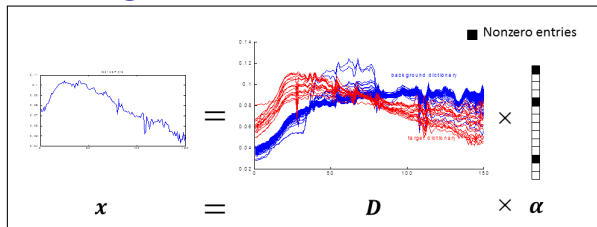


# Sparsity-based target detection<sup>6</sup>



<sup>6</sup>Chen et al., IEEE J. Sel. Topics Signal Process., 2011

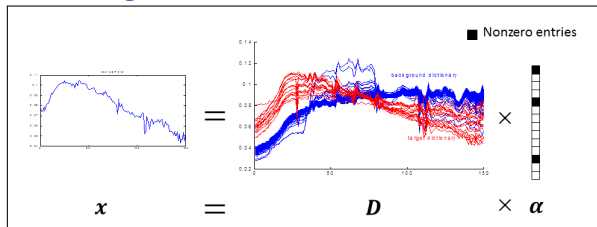
# Sparsity-based target detection<sup>6</sup>



$$x \approx D_t \alpha_t + D_b \alpha_b = \underbrace{\begin{bmatrix} D_t & D_b \end{bmatrix}}_D \begin{bmatrix} \alpha_t \\ \alpha_b \end{bmatrix} = D \alpha$$

- $D_t$  : target dictionary (matrix with columns as target spectra)
- $D_b$  : background dictionary (columns  $\rightarrow$  background spectra)

# Sparsity-based target detection<sup>6</sup>

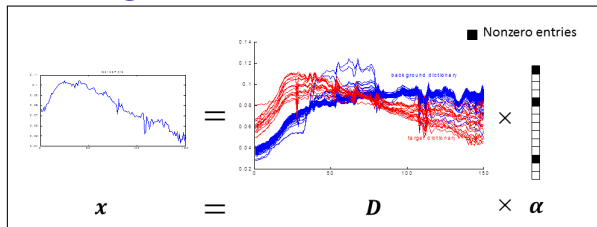


$$\mathbf{x} \approx \mathbf{D}_t \boldsymbol{\alpha}_t + \mathbf{D}_b \boldsymbol{\alpha}_b = \underbrace{\begin{bmatrix} \mathbf{D}_t & \mathbf{D}_b \end{bmatrix}}_{\mathbf{D}} \begin{bmatrix} \boldsymbol{\alpha}_t \\ \boldsymbol{\alpha}_b \end{bmatrix} = \mathbf{D} \boldsymbol{\alpha}$$

- $\mathbf{D}_t$  : target dictionary (matrix with columns as target spectra)
- $\mathbf{D}_b$  : background dictionary (columns  $\rightarrow$  background spectra)

**Sparse recovery:**  $\hat{\boldsymbol{\alpha}} = \arg \min \|\boldsymbol{\alpha}\|_0$  subject to  $\|\mathbf{x} - \mathbf{D}\boldsymbol{\alpha}\|_2 \leq \epsilon$

# Sparsity-based target detection<sup>6</sup>



$$\mathbf{x} \approx \mathbf{D}_t \boldsymbol{\alpha}_t + \mathbf{D}_b \boldsymbol{\alpha}_b = \underbrace{\begin{bmatrix} \mathbf{D}_t & \mathbf{D}_b \end{bmatrix}}_{\mathbf{D}} \begin{bmatrix} \boldsymbol{\alpha}_t \\ \boldsymbol{\alpha}_b \end{bmatrix} = \mathbf{D} \boldsymbol{\alpha}$$

- $\mathbf{D}_t$ : target dictionary (matrix with columns as target spectra)
- $\mathbf{D}_b$ : background dictionary (columns  $\rightarrow$  background spectra)

**Sparse recovery:**  $\hat{\boldsymbol{\alpha}} = \arg \min \|\boldsymbol{\alpha}\|_0$  subject to  $\|\mathbf{x} - \mathbf{D}\boldsymbol{\alpha}\|_2 \leq \epsilon$

- **Class assignment:**

$$R(\mathbf{x}) = \|\mathbf{x} - \mathbf{D}_b \hat{\boldsymbol{\alpha}}_b\|_2 - \|\mathbf{x} - \mathbf{D}_t \hat{\boldsymbol{\alpha}}_t\|_2$$

$$\mathbf{x} = \begin{cases} \text{target,} & \text{if } R(\mathbf{x}) > \delta \\ \text{background,} & \text{otherwise.} \end{cases}$$

## Joint sparsity model for HSI target detection<sup>7,8</sup>

- Homogeneous regions in HSI → neighboring pixels strongly correlated → **same** sparsity pattern, weighted **differently**

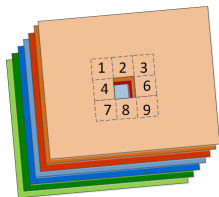
---

<sup>7</sup> Chen et al., IEEE Geosci. Remote Sens. Lett., 2011

<sup>8</sup> Chen et al., IEEE Trans. Geosci. Remote Sens., 2011

# Joint sparsity model for HSI target detection<sup>7,8</sup>

- Homogeneous regions in HSI → neighboring pixels strongly correlated → **same** sparsity pattern, weighted **differently**



$$\mathbf{x}_i = D\boldsymbol{\alpha}_i, i = 1, 2, \dots, T$$

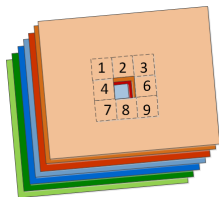
$$\begin{aligned} \mathbf{X} &= [\mathbf{x}_1 \quad \mathbf{x}_2 \quad \cdots \quad \mathbf{x}_T] \\ &= D \underbrace{[\boldsymbol{\alpha}_1 \quad \boldsymbol{\alpha}_2 \quad \cdots \quad \boldsymbol{\alpha}_T]}_S = DS \end{aligned}$$

<sup>7</sup> Chen et al., IEEE Geosci. Remote Sens. Lett., 2011

<sup>8</sup> Chen et al., IEEE Trans. Geosci. Remote Sens., 2011

# Joint sparsity model for HSI target detection<sup>7,8</sup>

- Homogeneous regions in HSI  $\rightarrow$  neighboring pixels strongly correlated  $\rightarrow$  same sparsity pattern, weighted differently



$$\mathbf{x}_i = D\boldsymbol{\alpha}_i, i = 1, 2, \dots, T$$

$$\begin{aligned} \mathbf{X} &= [\mathbf{x}_1 \quad \mathbf{x}_2 \quad \cdots \quad \mathbf{x}_T] \\ &= D \underbrace{[\boldsymbol{\alpha}_1 \quad \boldsymbol{\alpha}_2 \quad \cdots \quad \boldsymbol{\alpha}_T]}_{\mathbf{S}} = D\mathbf{S} \end{aligned}$$

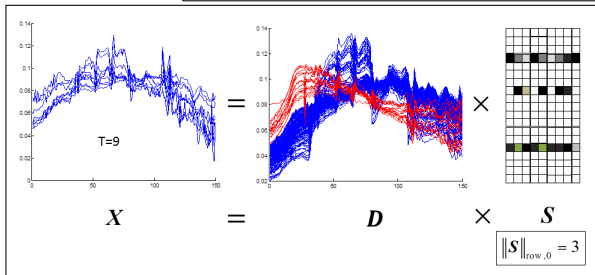


Figure :  $\mathbf{S}$  has only a few nonzero rows.

<sup>7</sup> Chen et al., IEEE Geosci. Remote Sens. Lett., 2011

<sup>8</sup> Chen et al., IEEE Trans. Geosci. Remote Sens., 2011

# Joint sparse recovery

$$\hat{\mathbf{S}} = \arg \min \|\mathbf{DS} - \mathbf{X}\|_F \quad \text{subject to} \quad \|\mathbf{S}\|_{\text{row},0} \leq K_0$$

- $\|\mathbf{S}\|_{\text{row},0} \rightarrow$  number of nonzero rows in  $\mathbf{S}$

---

<sup>9</sup>Tropp et al., Signal Processing, 2006

<sup>10</sup>Tropp, Signal Processing, 2006



# Joint sparse recovery

$$\hat{\mathbf{S}} = \arg \min \|\mathbf{D}\mathbf{S} - \mathbf{X}\|_F \quad \text{subject to} \quad \|\mathbf{S}\|_{\text{row},0} \leq K_0$$

- $\|\mathbf{S}\|_{\text{row},0} \rightarrow$  number of nonzero rows in  $\mathbf{S}$
- Recovery algorithms:
  - Simultaneous versions of greedy pursuit algorithms (SOMP) <sup>9</sup>
  - Convex relaxation of  $\|\mathbf{S}\|_{\text{row},0}$  to mixed  $\ell_1/\ell_2$ -norm<sup>10</sup>

---

<sup>9</sup>Tropp et al., Signal Processing, 2006

<sup>10</sup>Tropp, Signal Processing, 2006

# Joint sparse recovery

$$\hat{\mathbf{S}} = \arg \min \|\mathbf{D}\mathbf{S} - \mathbf{X}\|_F \quad \text{subject to} \quad \|\mathbf{S}\|_{\text{row},0} \leq K_0$$

- $\|\mathbf{S}\|_{\text{row},0} \rightarrow$  number of nonzero rows in  $\mathbf{S}$
- Recovery algorithms:
  - Simultaneous versions of greedy pursuit algorithms (SOMP) <sup>9</sup>
  - Convex relaxation of  $\|\mathbf{S}\|_{\text{row},0}$  to mixed  $\ell_1/\ell_2$ -norm<sup>10</sup>
- Class decision using total reconstruction residual:

$$r(\mathbf{x}) = \|\mathbf{X} - \mathbf{D}_b \hat{\mathbf{S}}_b\|_F - \|\mathbf{X} - \mathbf{D}_t \hat{\mathbf{S}}_t\|_F$$

---

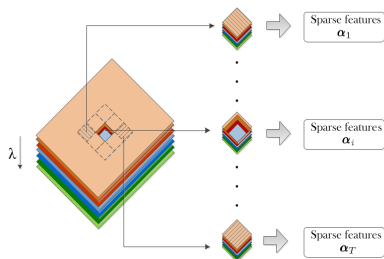
<sup>9</sup>Tropp et al., Signal Processing, 2006

<sup>10</sup>Tropp, Signal Processing, 2006

# Summary: Opportunities and challenges

- 1 Feature representations of pixels in a local neighborhood are statistically correlated
  - How to mine the **class-conditional correlations** among distinct feature representations?
- 2 Practical concern: Availability of very few target pixels
  - How to learn class-conditional models under **limited training**?
- 3 Use of reconstruction residual for detection
  - Design of truly discriminative classifiers

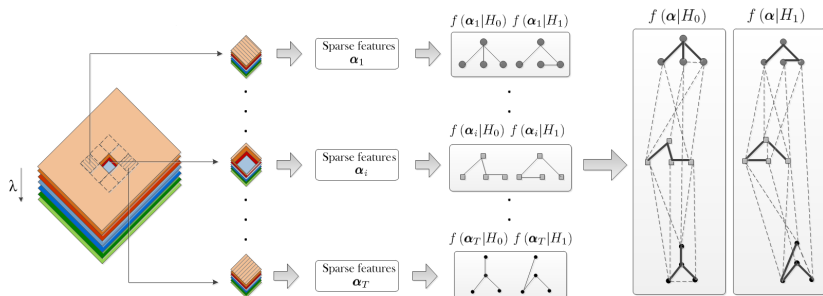
# Spatio-spectral sparsity via discriminative graphical models



Two-stage framework:

- 1 Extract multiple local sparse feature representations, which are **conditionally correlated** per class

# Spatio-spectral sparsity via discriminative graphical models



Two-stage framework:

- 1 Extract multiple local sparse feature representations, which are **conditionally correlated** per class
- 2 Mine dependencies between different features via boosting on discriminative graphs

## Probabilistic graphical models: A brief review

- **Graph**  $\mathcal{G} = (\mathcal{V}, \mathcal{E})$  defined by a set of nodes  $\mathcal{V} = \{1, \dots, n\}$ , and a set of edges  $\mathcal{E} \subset \binom{\mathcal{V}}{2}$ .
- **Graphical model**: Random vector defined on a graph; nodes represent random variables, edges reveal conditional dependencies.
- Graph structure defines factorization of joint probability distribution

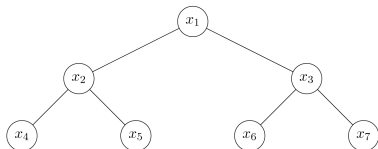


Figure : Tree - connected acyclic graph.

$$f(\mathbf{x}) = f(x_1)f(x_2|x_1)f(x_3|x_1)f(x_4|x_2)f(x_5|x_2)f(x_6|x_3)f(x_7|x_3).$$



$f(\alpha_i H_0), f(\alpha_i H_1)$
------------------------------------

# Learning graphical models

- **Generative learning**: Single graph to minimize **approximation** error<sup>11</sup>

Given  $p$ , find  $\hat{p} = \arg \min_{\hat{p} \text{ is a tree}} D(p||\hat{p})$ .

$$\left( D(p||\hat{p}) := \int p(\mathbf{x}) \log \left( \frac{p(\mathbf{x})}{\hat{p}(\mathbf{x})} \right) d\mathbf{x} \rightarrow \text{KL-divergence.} \right)$$

- **Discriminative learning**: Simultaneously learn a **pair** of graphs to approximately minimize **classification** error<sup>12</sup>

---

<sup>11</sup>Chow et al., IEEE Trans. Inf. Theory, 1968

<sup>12</sup>Tan et al., IEEE Trans. Signal Process., 2010

# Learning graphical models

- **Generative learning:** Single graph to minimize **approximation** error<sup>11</sup>

Given  $p$ , find  $\hat{p} = \arg \min_{\hat{p} \text{ is a tree}} D(p||\hat{p})$ .

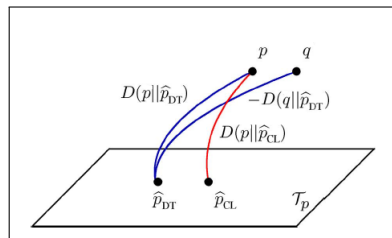
$$\left( D(p||\hat{p}) := \int p(\mathbf{x}) \log \left( \frac{p(\mathbf{x})}{\hat{p}(\mathbf{x})} \right) d\mathbf{x} \rightarrow \text{KL-divergence.} \right)$$

- **Discriminative learning:** Simultaneously learn a **pair** of graphs to approximately minimize **classification** error<sup>12</sup>

Tree-approximate  $J$ -divergence:

$$\hat{J}(\hat{p}, \hat{q}; p, q) := \int_{\Omega \subset \mathcal{X}^n} (p(\mathbf{x}) - q(\mathbf{x})) \log \left( \frac{\hat{p}(\mathbf{x})}{\hat{q}(\mathbf{x})} \right) d\mathbf{x}.$$

$$(\hat{p}, \hat{q}) = \arg \max_{\hat{p} \in \mathcal{T}_p, \hat{q} \in \mathcal{T}_q} \hat{J}(\hat{p}, \hat{q}; p, q).$$



(Figure courtesy Tan et al.)

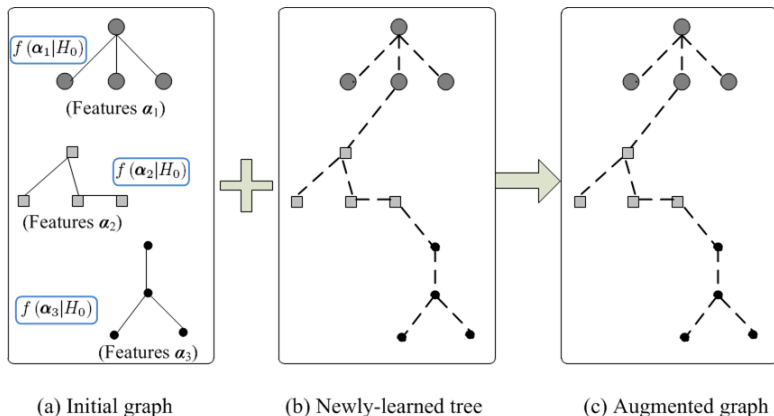
<sup>11</sup>Chow et al., IEEE Trans. Inf. Theory, 1968

<sup>12</sup>Tan et al., IEEE Trans. Signal Process., 2010



# Learning discriminative graphs: An illustration<sup>13</sup>

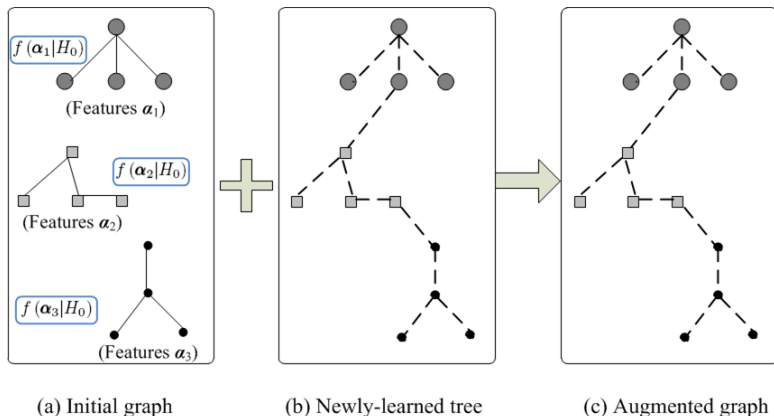
Iteration 1:



<sup>13</sup> Shown for distribution  $p$ ; graph for  $q$  learnt analogously.

# Learning discriminative graphs: An illustration<sup>13</sup>

Iteration 1:

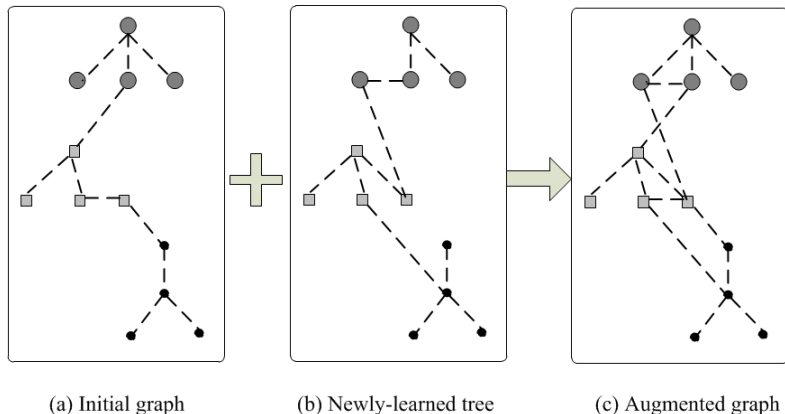


Re-weighting of training samples (boosting)  $\rightarrow$  learn another tree ...

<sup>13</sup> Shown for distribution  $p$ ; graph for  $q$  learnt analogously.

# Learning discriminative graphs: An illustration

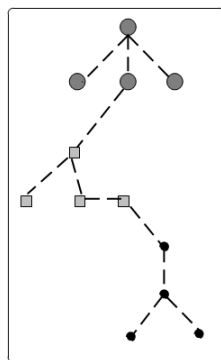
Iteration 2:



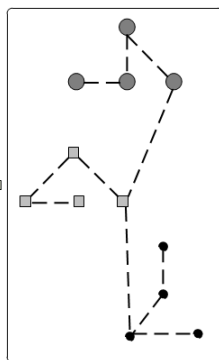
Newly introduced edges crucial for capturing correlations amongst distinct signal representations.

# Learning discriminative graphs: An illustration

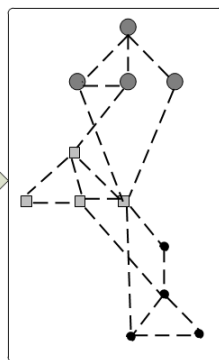
Iteration 3:



(a) Initial graph



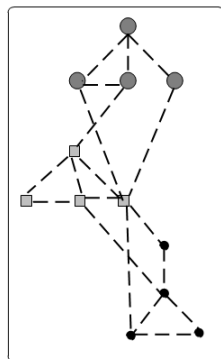
(b) Newly-learned tree



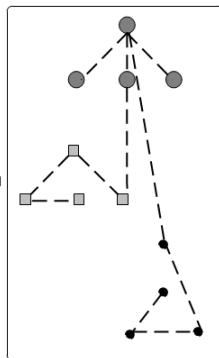
(c) Augmented graph

# Learning discriminative graphs: An illustration

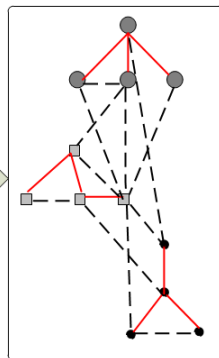
Iteration 4:



(a) Initial graph



(b) Newly-learned tree



(c) Augmented graph

# Stopping criterion and class assignment

How many edges to learn?

- 1 Cross-validation
- 2 Using the  $J$ -divergence:

$$\hat{J}(\hat{p}, \hat{q}; p, q) := \int_{\Omega \subset \mathcal{X}^n} (p(\mathbf{x}) - q(\mathbf{x})) \log \left( \frac{\hat{p}(\mathbf{x})}{\hat{q}(\mathbf{x})} \right) d\mathbf{x}.$$

Stopping criterion:

Stop after  $i$  boosting iterations if:

$$\frac{\hat{J}^{(i+1)}(\hat{p}, \hat{q}; p, q) - \hat{J}^{(i)}(\hat{p}, \hat{q}; p, q)}{\hat{J}^{(i)}(\hat{p}, \hat{q}; p, q)} < \epsilon$$

# Stopping criterion and class assignment

How many edges to learn?

- 1 Cross-validation
- 2 Using the  $J$ -divergence:

$$\hat{J}(\hat{p}, \hat{q}; p, q) := \int_{\Omega \subset \mathcal{X}^n} (p(\mathbf{x}) - q(\mathbf{x})) \log \left( \frac{\hat{p}(\mathbf{x})}{\hat{q}(\mathbf{x})} \right) d\mathbf{x}.$$

Stopping criterion:

Stop after  $i$  boosting iterations if:

$$\frac{\hat{J}^{(i+1)}(\hat{p}, \hat{q}; p, q) - \hat{J}^{(i)}(\hat{p}, \hat{q}; p, q)}{\hat{J}^{(i)}(\hat{p}, \hat{q}; p, q)} < \epsilon$$

Class label assignment:

$$\text{Class}(\mathbf{x}) = \begin{cases} \text{Target} & \text{if } \log \left( \frac{\hat{f}(\boldsymbol{\alpha}|H_1)}{\hat{f}(\boldsymbol{\alpha}|H_0)} \right) \geq \tau \\ \text{Background} & \text{if } \log \left( \frac{\hat{f}(\boldsymbol{\alpha}|H_1)}{\hat{f}(\boldsymbol{\alpha}|H_0)} \right) < \tau. \end{cases}$$

# Experiment I

- HYDICE desert radiance II (DR-II) data collection<sup>14</sup>



- 150 spectral bands: 400nm - 2500nm; 6 targets
- Target dictionary  $D_t$ : 18 training spectra from leftmost target
- Background dictionary  $D_b$ : 216 training spectra

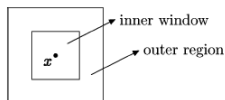


Figure : Dual window for adaptive background selection<sup>15</sup>

<sup>14</sup>Basedow et al., SPIE Conf. Algorithms Technol. Multispectral, Hyperspectral, Ultraspectral Imagery XV, 1995

<sup>15</sup>Chen et al., IEEE J. Sel. Topics Signal Process., 2011



# Results: Confusion matrix

Methods compared:

- 1 **MSD**: matched subspace detector<sup>16</sup>
- 2 **SVM-CK**: composite kernel SVM<sup>17</sup>
- 3 **SOMP**: joint sparsity model<sup>18</sup>
- 4 **LSGM**: proposed local-sparsity-graphical-model approach

Table : Confusion matrix for DR-II hyperspectral image.

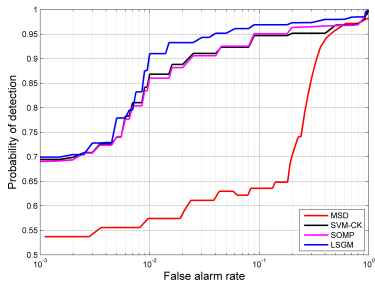
Class	Target	Background	Method
Target	0.6356	0.3644	MSD
	0.9469	0.0531	SVM-CK
	0.9508	0.0492	SOMP
	<b>0.9690</b>	<b>0.0310</b>	<b>LSGM</b>
Background	0.0157	0.9843	MSD
	0.0075	0.9925	SVM-CK
	0.0077	0.9923	SOMP
	<b>0.0074</b>	<b>0.9926</b>	<b>LSGM</b>

<sup>16</sup>Scharf et al., IEEE Trans. Signal Process., 1994

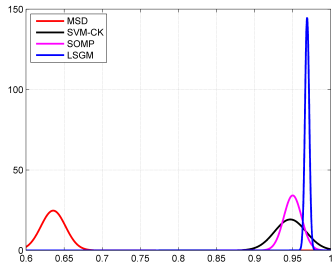
<sup>17</sup>Camps-Valls et al., IEEE Geosci. Remote Sens. Lett, 2006

<sup>18</sup>Chen et al., IEEE Trans. Geosci. Remote Sens., 2011

# Results: Receiver operating characteristic (ROC)



(a) ROC for DR-II.



(b) Density function of detection rates (multiple training runs).

## Experiment II

- HYDICE forest radiance I (FR-I) data collection

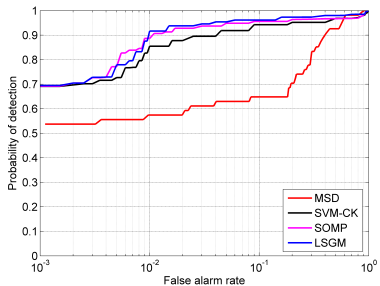


- 150 spectral bands: 400nm - 2500nm; 14 targets

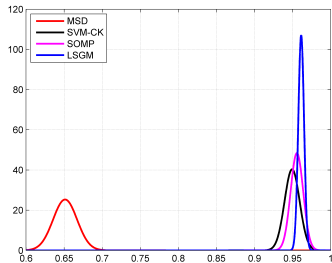
Table : Confusion matrix for FR-I hyperspectral image.

Class	Target	Background	Method
Target	0.6512	0.3488	MSD
	0.9493	0.0507	SVM-CK
	0.9556	0.0444	SOMP
	<b>0.9612</b>	<b>0.0388</b>	<b>LSGM</b>
Background	0.0239	0.9761	MSD
	0.0090	0.9910	SVM-CK
	0.0097	0.9903	SOMP
	<b>0.0086</b>	<b>0.9914</b>	<b>LSGM</b>

# Results: ROC



(a) ROC for FR-I.



(b) Density function of detection rates (multiple training runs).

# Conclusions

- Probabilistic graphical model framework for hyperspectral target detection
  - Incorporates **spatio-spectral** notion of sparsity via joint sparsity model
  - Explicitly captures **conditional correlations** between local sparse features (instead of using reconstruction residuals) for better discrimination.

\* Journal version (HSI classification) accepted to IEEE Geosci. Remote Sens. Lett., June 2012

Thank you

Questions?

## Backup Slides

Edge weights:

$$\begin{aligned}\psi_{i,j}^p &:= \mathbb{E}_{\tilde{p}_{i,j}} \left[ \log \frac{\tilde{p}_{i,j}}{\tilde{p}_i \tilde{p}_j} \right] - \mathbb{E}_{\tilde{q}_{i,j}} \left[ \log \frac{\tilde{p}_{i,j}}{\tilde{p}_i \tilde{p}_j} \right] \\ \psi_{i,j}^q &:= \mathbb{E}_{\tilde{q}_{i,j}} \left[ \log \frac{\tilde{q}_{i,j}}{\tilde{q}_i \tilde{q}_j} \right] - \mathbb{E}_{\tilde{p}_{i,j}} \left[ \log \frac{\tilde{q}_{i,j}}{\tilde{q}_i \tilde{q}_j} \right].\end{aligned}$$

---

**Algorithm 1** Discriminative trees (DT)

---

Given: Training sets  $\mathcal{T}_p$  and  $\mathcal{T}_q$ .

- 1: Estimate pairwise statistics  $\tilde{p}_{i,j}(x_i, x_j)$ ,  $\tilde{q}_{i,j}(x_i, x_j)$  for all edges  $(i, j)$ .
  - 2: Compute edge weights  $\psi_{i,j}^p$  and  $\psi_{i,j}^q$  for all edges  $(i, j)$ .
  - 3: Find  $\mathcal{E}_{\hat{p}} = \text{MWST}(\psi_{i,j}^p)$  and  $\mathcal{E}_{\hat{q}} = \text{MWST}(\psi_{i,j}^q)$ .
  - 4: Get  $\hat{p}$  by projection of  $\tilde{p}$  onto  $\mathcal{E}_{\hat{p}}$ ; likewise  $\hat{q}$ .
  - 5: LRT using  $\hat{p}$  and  $\hat{q}$ .
-



---

## Algorithm 2 AdaBoost learning algorithm

---

- 1: Input data  $(x_i, y_i)$ ,  $i = 1, 2, \dots, N$ , where  $x_i \in S$ ,  $y_i \in \{-1, +1\}$
  - 2: Initialize  $D_1(i) = \frac{1}{N}$ ,  $i = 1, 2, \dots, N$
  - 3: For  $t = 1, 2, \dots, T$ :
    - Train weak learner using distribution  $D_t$
    - Determine weak hypothesis  $h_t : S \mapsto \mathbb{R}$  with error  $\epsilon_t$
    - Choose  $\beta_t = \frac{1}{2} \ln \left( \frac{1-\epsilon_t}{\epsilon_t} \right)$
    - $D_{t+1}(i) = \frac{1}{Z_t} \{D_t(i) \exp(-\beta_t y_i h_t(x_i))\}$ , where  $Z_t$  is a normalization factor
  - 4: Output soft decision  $H(x) = \text{sign} \left[ \sum_{t=1}^T \beta_t h_t(x) \right]$ .
- 

- Distribution of weights over the training set
- In each iteration, weak learner  $h_t$  minimizes weighted training error
- Weights on incorrectly classified samples increased  $\rightarrow$  slow learners penalized for harder examples.

# Learning thicker graphical models

- Final boosted classifier:

$$\begin{aligned} H_T(\mathbf{x}) &= \operatorname{sgn} \left[ \sum_{t=1}^T \alpha_t \log \left( \frac{\hat{p}_t(\mathbf{x})}{\hat{q}_t(\mathbf{x})} \right) \right] = \operatorname{sgn} \left[ \log \prod_{t=1}^T \left( \frac{\hat{p}_t(\mathbf{x})}{\hat{q}_t(\mathbf{x})} \right)^{\alpha_t} \right] \\ &= \operatorname{sgn} \left[ \log \left( \frac{\prod_{t=1}^T (\hat{p}_t(\mathbf{x}))^{\alpha_t}}{\prod_{t=1}^T (\hat{q}_t(\mathbf{x}))^{\alpha_t}} \right) \right] = \operatorname{sgn} \left[ \log \left( \frac{\hat{p}(\mathbf{x})}{\hat{q}(\mathbf{x})} \right) \right] \end{aligned}$$

Define:

$$Z_p(\boldsymbol{\alpha}) = Z_p(\alpha_1, \dots, \alpha_T) = \sum_{\mathbf{x}} \hat{p}(\mathbf{x}); Z_q(\boldsymbol{\alpha}) = \sum_{\mathbf{x}} \hat{q}(\mathbf{x})$$

- Normalized distributions for inference:  $\frac{\hat{p}(\mathbf{x})}{Z_p(\boldsymbol{\alpha})}, \frac{\hat{q}(\mathbf{x})}{Z_q(\boldsymbol{\alpha})}$

→ Thicker **graphical models** learnt.

# Local Sparsity Graphical Models

---

## Algorithm 3 LSGM (Steps 1-4 offline)

---

- 1: **Feature extraction (training):** Compute sparse representations  $\alpha_l, l = 1, \dots, T$  for neighboring pixels of the training data
  - 2: **Initial disjoint graphs:**  
Discriminatively learn  $T$  pairs of  $N$ -node tree graphs  $\mathcal{G}_l^t$  and  $\mathcal{G}_l^b$  on  $\{\alpha_l\}$ , for  $l = 1, \dots, T$ , obtained from training data
  - 3: Separately concatenate nodes corresponding to the two classes, to generate initial graphs
  - 4: **Boosting on disjoint graphs:** Iteratively thicken initial disjoint graphs via boosting to obtain final graphs  $\mathcal{G}^t$  and  $\mathcal{G}^b$
- 
- {**Online process**}
- 5: **Feature extraction (test):** Obtain sparse representations  $\alpha_l, l = 1, \dots, T$  in  $\mathbb{R}^N$  from test image
  - 6: **Inference:** Classify based on output of the resulting classifier.
-

# Simultaneous Orthogonal Matching Pursuit

---

## Algorithm 4 SOMP

---

**Input:**  $B \times N$  dictionary matrix  $\mathbf{D} = [\mathbf{d}_1 \cdots \mathbf{d}_N]$ ,  $B \times T$  signal matrix  $\mathbf{X} = [\mathbf{x}_1 \cdots \mathbf{x}_T]$ , and number of iterations  $K$

Initialization: residual  $\mathbf{R}_0 = \mathbf{X}$ , index set  $\Lambda_0 = \emptyset$ , iteration counter  $k = 1$

**while**  $k \leq K$  **do**

(1) Find the index of the atom that best approximates all residuals:  $\lambda_k = \arg \max_{i=1, \dots, N} \|\mathbf{R}_{k-1}^T \mathbf{d}_i\|_p, p \geq 1$

(2) Update the index set  $\Lambda_k = \Lambda_{k-1} \cup \{\lambda_k\}$

(3) Compute the orthogonal projector  $\mathbf{P}_k = (\mathbf{D}_{\Lambda_k}^T \mathbf{D}_{\Lambda_k})^{-1} \mathbf{D}_{\Lambda_k}^T \mathbf{X} \in \mathbb{R}^{k \times T}$

where  $\mathbf{D}_{\Lambda_k} \in \mathbb{R}^{B \times k}$  consists of the  $k$  atoms in  $\mathbf{D}$  indexed in  $\Lambda_k$

(4) Update the residual matrix  $\mathbf{R}_k = \mathbf{X} - \mathbf{D}_{\Lambda_k} \mathbf{P}_k$

(5) Increment  $k$ :  $k \leftarrow k + 1$

**end while**

**Output:** Index set  $\Lambda = \Lambda_K$ , the sparse representation  $\mathbf{S}$  whose nonzero rows index by  $\Lambda$  are the  $K$  rows of the matrix  $(\mathbf{D}_{\Lambda_K}^T \mathbf{D}_{\Lambda_K})^{-1} \mathbf{D}_{\Lambda_K}^T \mathbf{X}$

---