

SENSING-AWARE CLASSIFICATION WITH HIGH DIMENSIONAL DATA

Prakash Ishwar

Burkay Orten, W. Clem Karl, Venkatesh Saligrama, Homer Pien



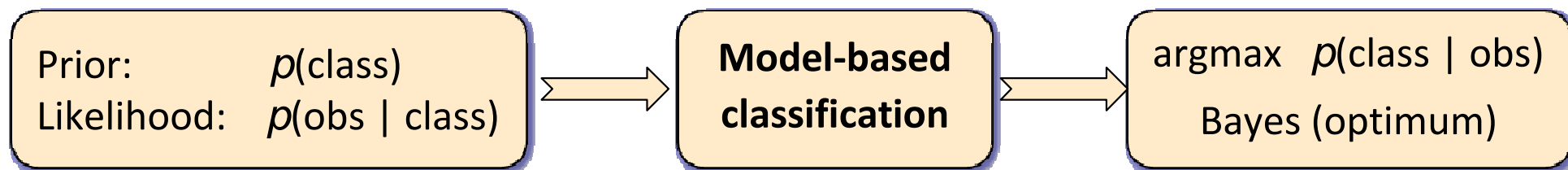
**BOSTON
UNIVERSITY**



Classification

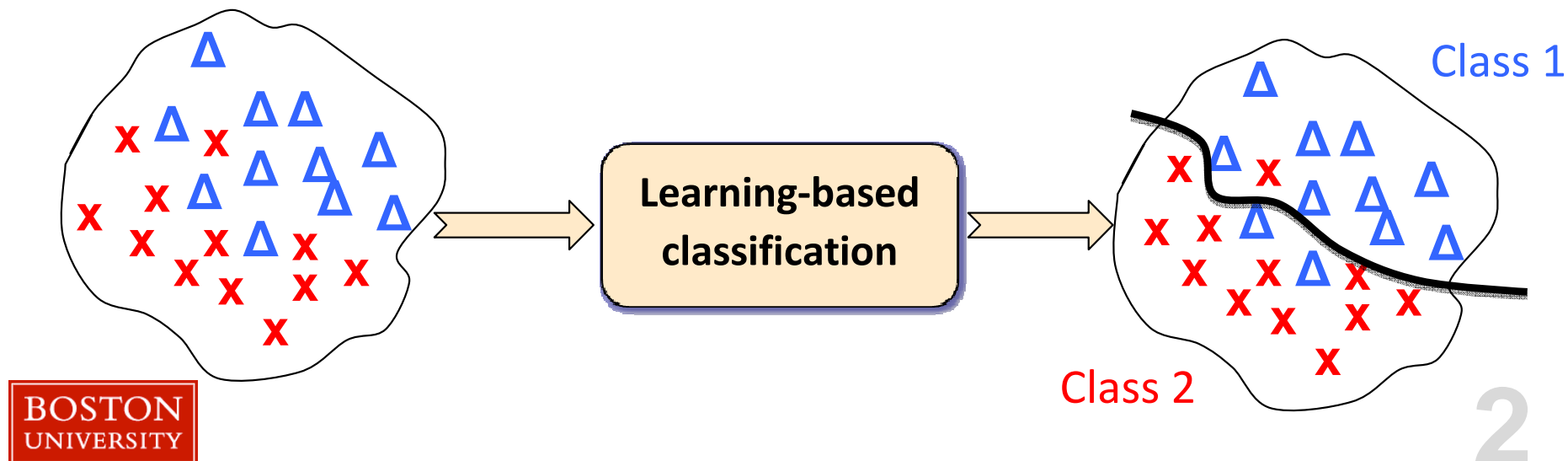
- Model-based classification**

Given: Model

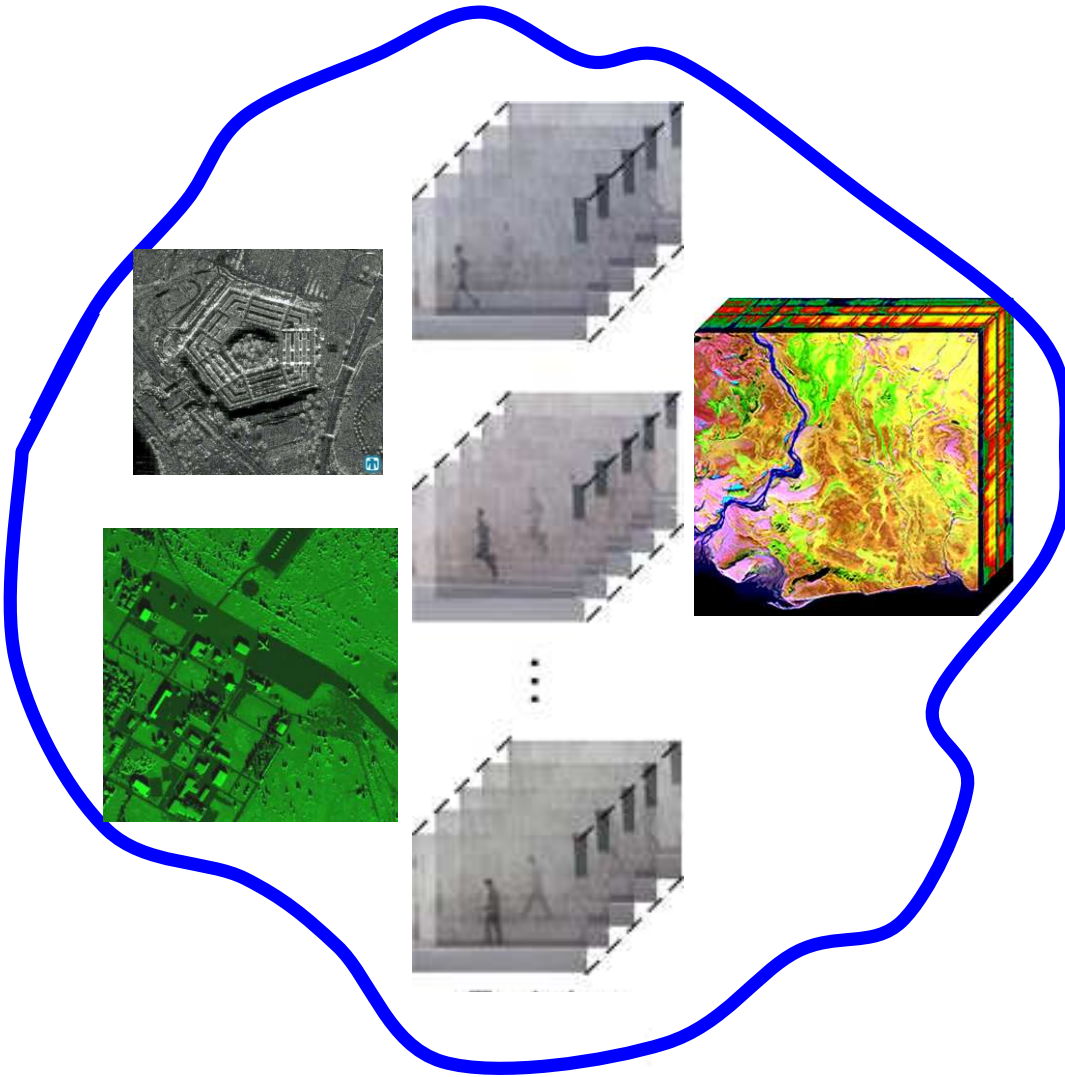


- Learning-based classification**

Given: Labeled Examples



Challenges and Approach



Challenges:

1. High dimensional data

- Surveillance video
- Hyperspectral images
- SAR
- MRI...

2. Few samples

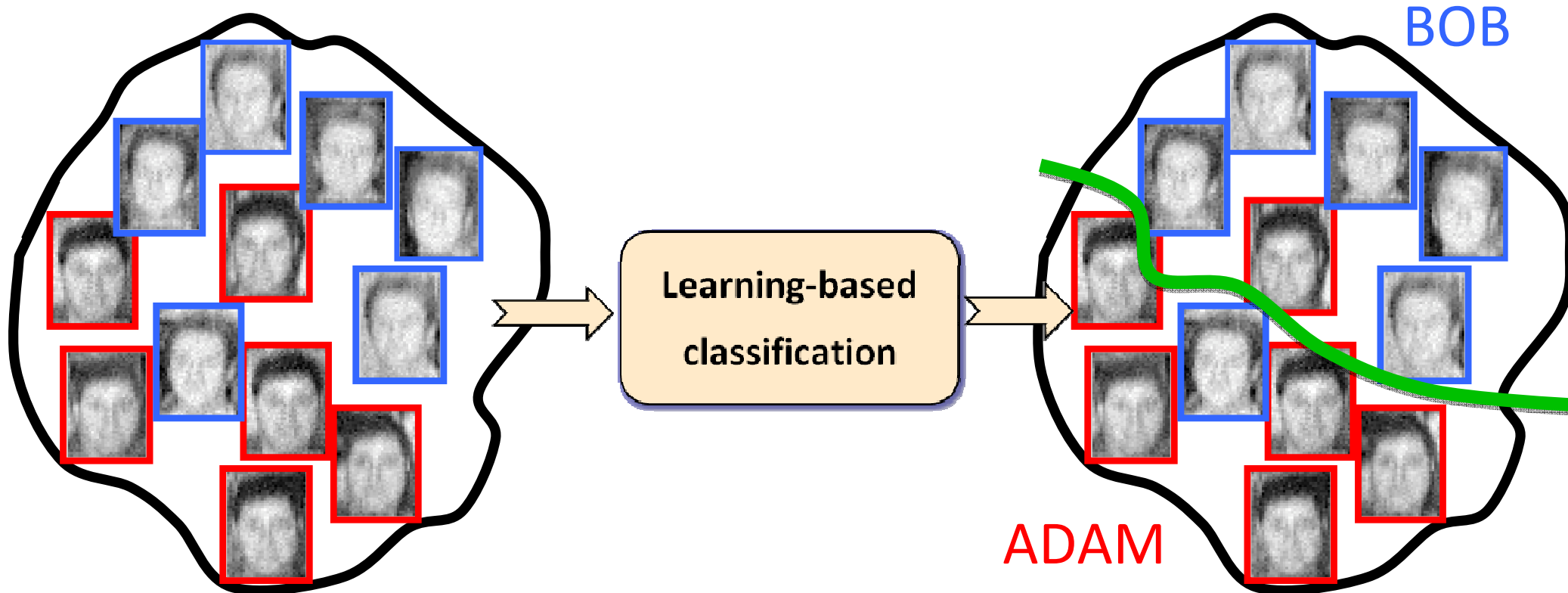
Approach:

- Exploit **latent low-dimensional sensing structure**

“Sensing-**Unaware**” Classification

Given: Blurry and noisy faces

Decision Rule



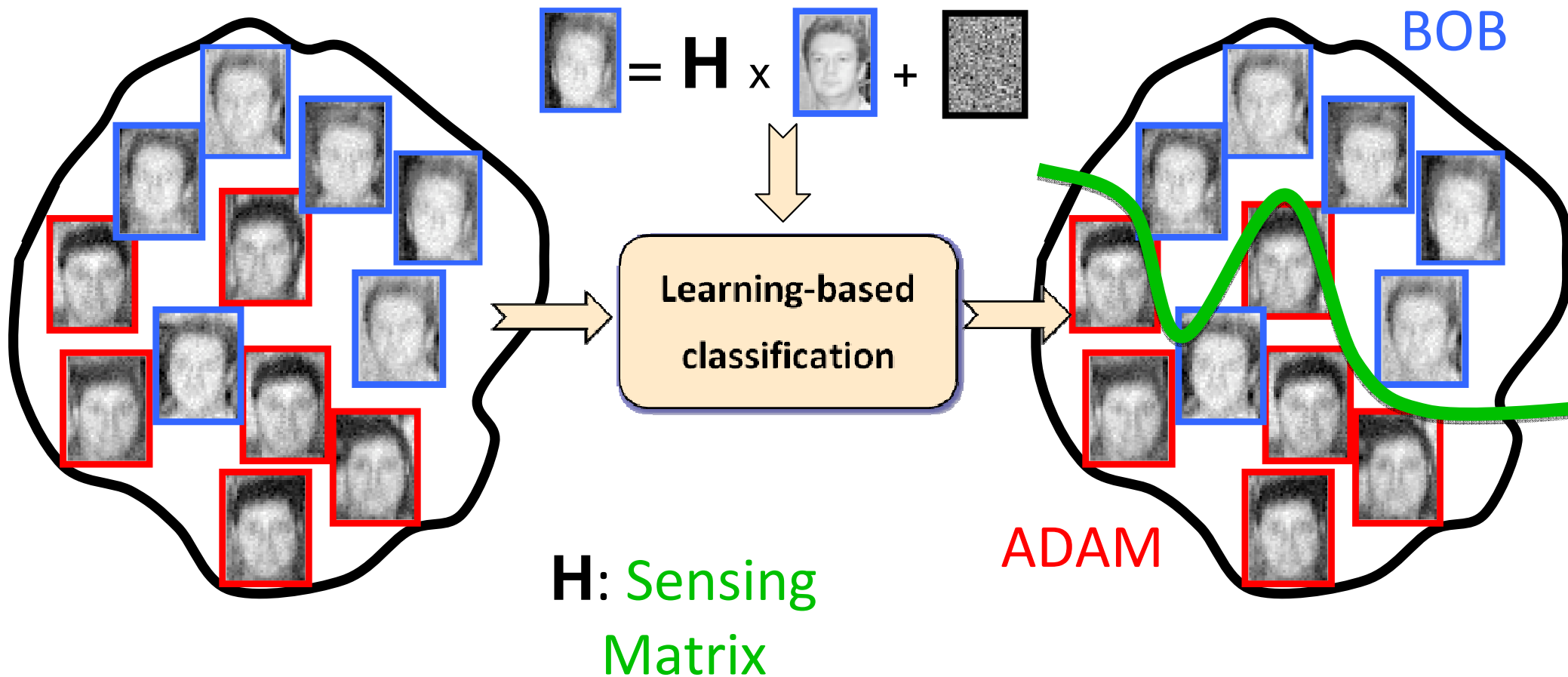
“Sensing-Aware” Classification

Given: Blurry and noisy faces

&

Sensing structure

Decision Rule



Asymptotic classification performance

- **Asymptotic regime of interest:**
 1. Data-dimension, # samples $\rightarrow \infty$
 2. Samples per dimension $\rightarrow 0$
 3. **Fixed problem difficulty:** Bayes risk kept fixed as dimension scales, i.e., problem not asymptotically “easy”.

- **Fundamental issue:** asymptotic classification performance
 - $P_{\text{error}} \rightarrow 1/2$ (random guessing)?
 - $P_{\text{error}} \rightarrow P_{\text{Bayes}}$?

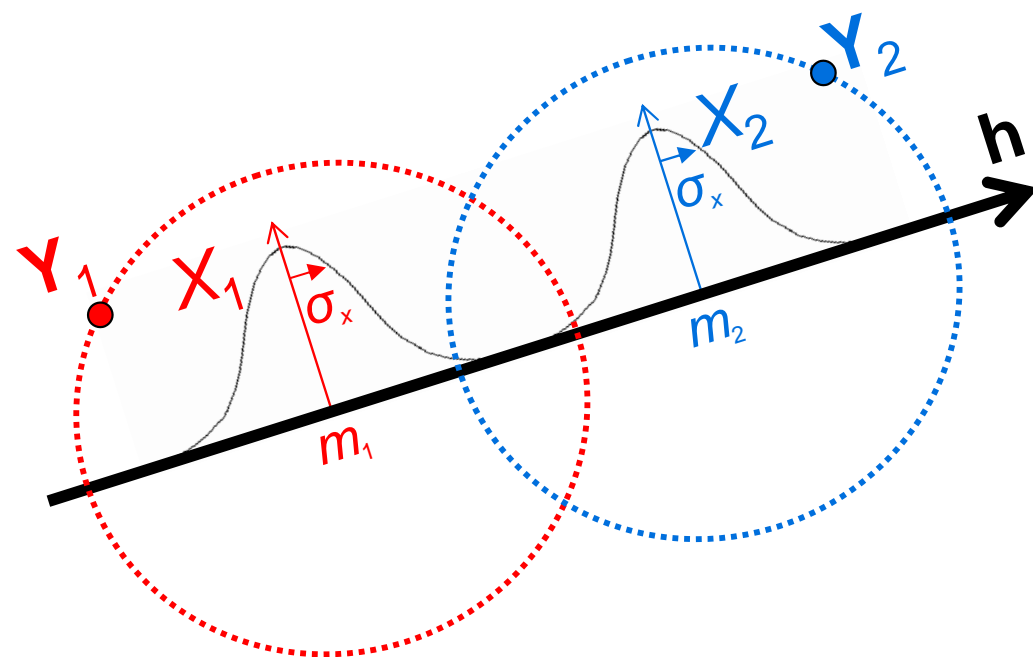
Mathematical model

- Sensing model:

$$\mathbf{Y}_i = \mathbf{h}X_i + \mathbf{Z}_i$$

$$\begin{bmatrix} \\ \\ \end{bmatrix}_{p \times 1} = \begin{bmatrix} \\ \\ \end{bmatrix}_{p \times 1} X_i + \begin{bmatrix} \\ \\ \end{bmatrix}_{p \times 1}$$

$i = 1, 2$ (class index)
equally likely



- $X_i \sim \mathcal{N}(m_i, \sigma_x^2), \quad \mathbf{Z}_i \sim \mathcal{N}(\mathbf{0}, \sigma_z^2 \mathbf{I})$

Mathematical model

- Sensing model:

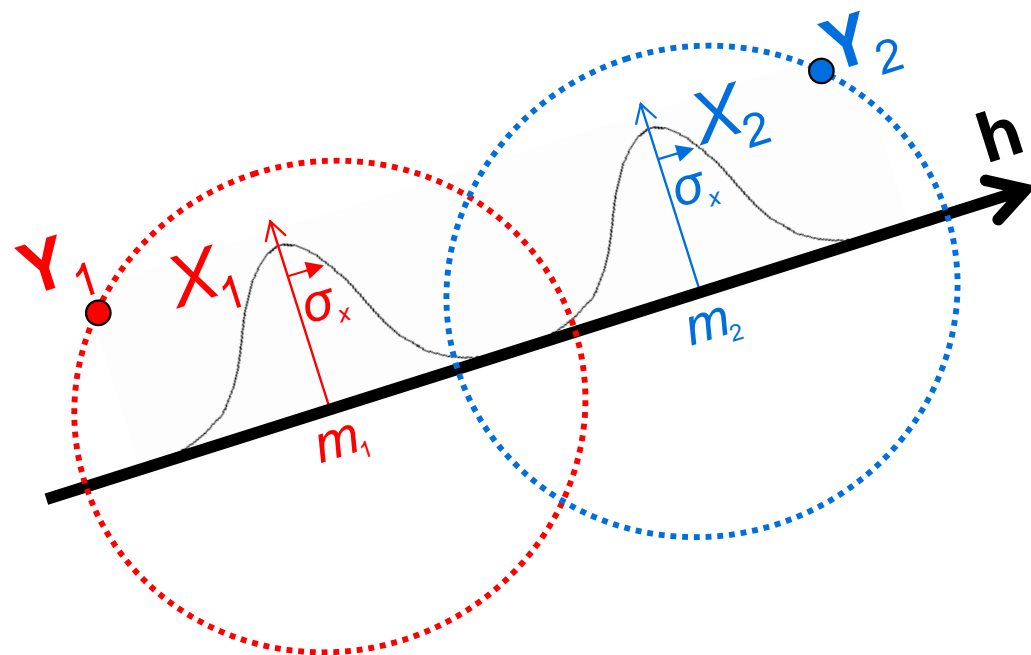
$$\mathbf{Y}_i = \mathbf{h}X_i + \mathbf{Z}_i$$

$$\begin{bmatrix} \\ \\ \end{bmatrix}_{p \times 1} = \begin{bmatrix} \\ \\ \end{bmatrix}_{p \times 1} X_i + \begin{bmatrix} \\ \\ \end{bmatrix}_{p \times 1}$$

$i = 1, 2$ (class index), equally likely

- Given: n labeled iid training examples, design classifier

$$\begin{pmatrix} \mathbf{Y}_{11}, \mathbf{Y}_{12}, \dots, \mathbf{Y}_{1n} \\ \mathbf{Y}_{21}, \mathbf{Y}_{22}, \dots, \mathbf{Y}_{2n} \end{pmatrix} \begin{matrix} \longleftarrow \text{Class 1 training examples} \\ \longleftarrow \text{Class 2 training examples} \end{matrix}$$

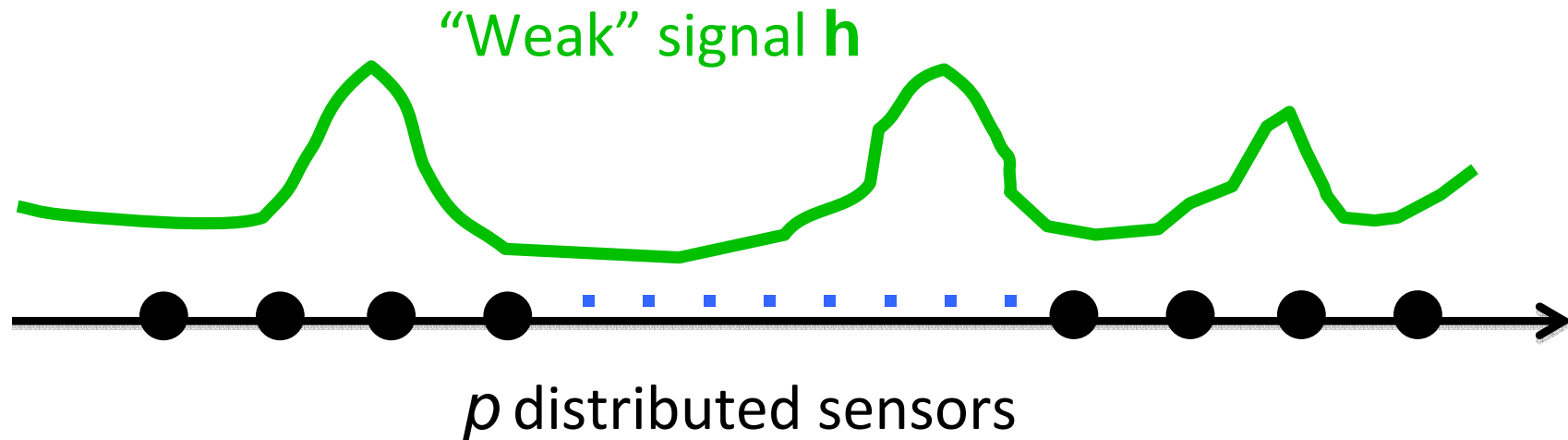


A sensor network application

Simple observation model using distributed sensors:

Class 1: $\mathbf{Y} = \mathbf{h} + \text{noise}$

Class 2: $\mathbf{Y} = -\mathbf{h} + \text{noise}$



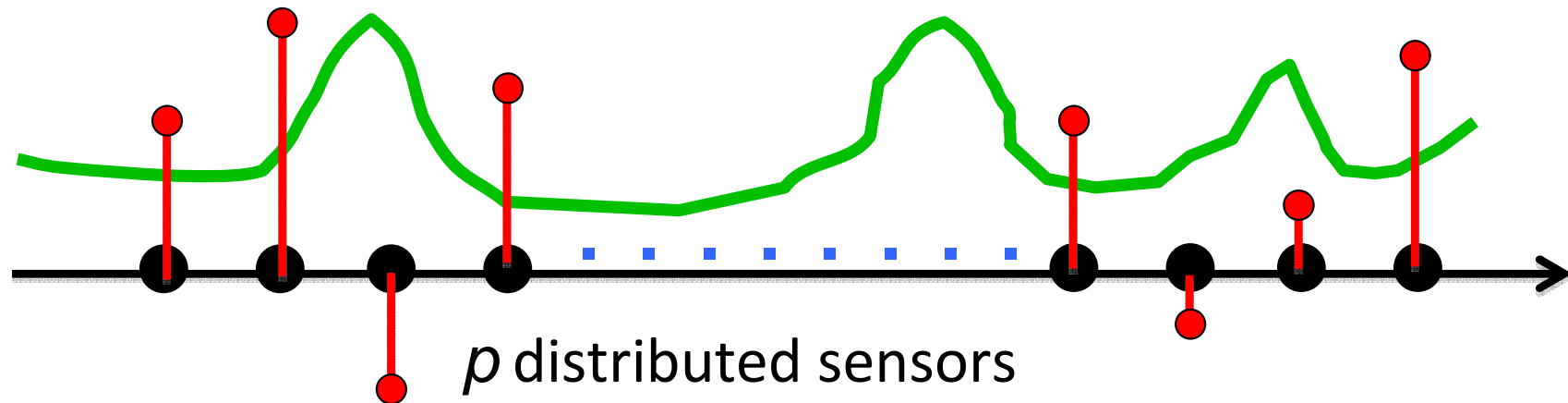
A sensor network application

Simple observation model using distributed sensors:

Class 1: $\mathbf{Y} = \mathbf{h} + \text{noise}$

Class 2: $\mathbf{Y} = -\mathbf{h} + \text{noise}$

Typical noisy observation given class 1

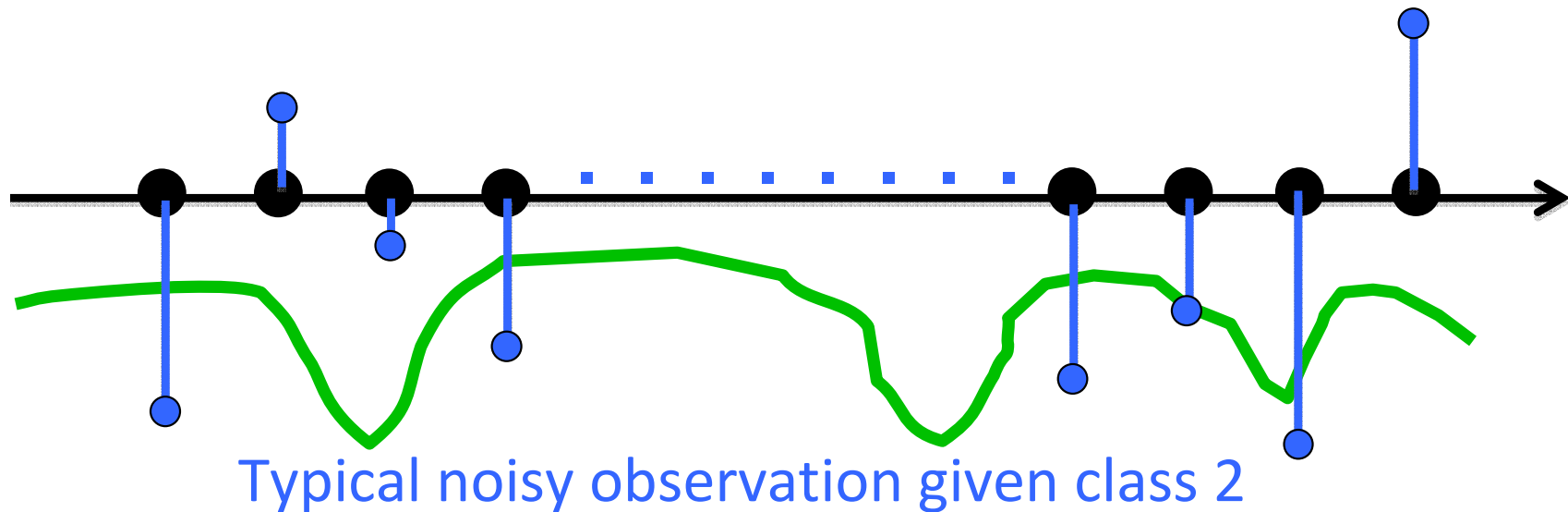


A sensor network application

Simple observation model using distributed sensors:

Class 1: $\mathbf{Y} = \mathbf{h} + \text{noise}$

Class 2: $\mathbf{Y} = -\mathbf{h} + \text{noise}$

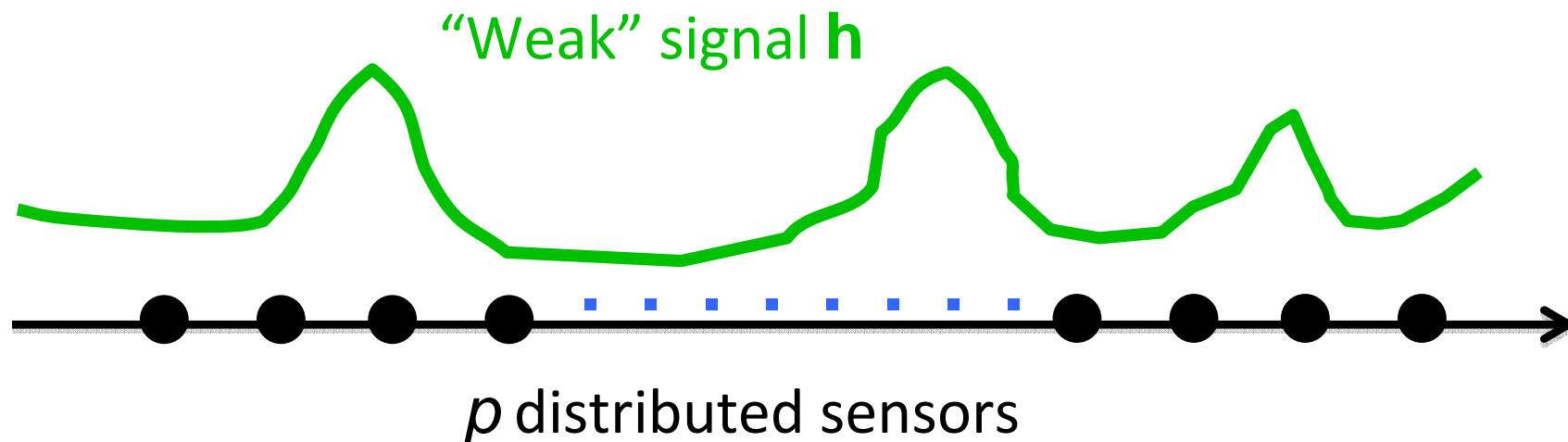


A sensor network application

Simple observation model using distributed sensors:

Class 1: $\mathbf{Y} = \mathbf{h} + \text{noise}$

Class 2: $\mathbf{Y} = -\mathbf{h} + \text{noise}$



Given:

n noisy observations of the "weak" signal per sensor from each class



Decide:

If a new set of observations are coming from class 1 or class 2

Classifiers Considered

- **Baseline:**
 1. Full Bayes
- **Unstructured:**
 2. Full ML
- **Structure-based:**
 3. Full structure
 4. Structured ML
 5. Structured sparsity

Sensing model

$$\mathbf{Y}_i = \mathbf{h}X_i + \mathbf{Z}_i \quad \begin{array}{l} i = 1 \text{ (class 1)} \\ i = 2 \text{ (class 2)} \end{array}$$

$$\begin{bmatrix} \\ \\ \end{bmatrix}_{p \times 1} = \begin{bmatrix} \\ \\ \end{bmatrix}_{p \times 1} X_i + \begin{bmatrix} \\ \\ \end{bmatrix}_{p \times 1}$$

Likelihood

$$p(\mathbf{Y} | \mathbf{h}, \text{class } i) = \mathcal{N}(\boldsymbol{\mu}_i, \boldsymbol{\Sigma})$$

1. Baseline: Full Bayes Classifier

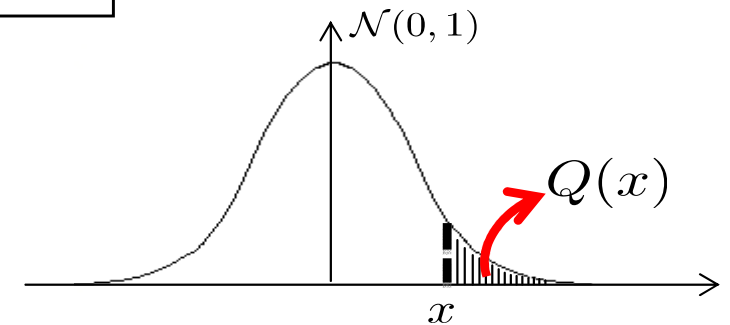
- Optimal test (Bayes rule): $\mathbf{1} \left\{ \frac{p(\mathbf{Y}|\mathbf{h}, \text{class } 2)}{p(\mathbf{Y}|\mathbf{h}, \text{class } 1)} \geq 1 \right\}$

- Classification rule: $\delta_{\text{Bayes}}(\mathbf{y}) = \mathbf{1} \left\{ \Delta^T \Sigma^{-1} (\mathbf{y} - \boldsymbol{\mu}) \geq 0 \right\}$

$$\begin{aligned} \Delta &= \mathbb{E}[\mathbf{Y}|\text{class } 2] - \mathbb{E}[\mathbf{Y}|\text{class } 1] &&= (m_2 - m_1)\mathbf{h} \\ \boldsymbol{\mu} &= \frac{1}{2} (\mathbb{E}[\mathbf{Y}|\text{class } 2] + \mathbb{E}[\mathbf{Y}|\text{class } 1]) &&= \frac{(m_2 + m_1)}{2} \mathbf{h} \\ \Sigma &= \text{cov}(\mathbf{Y}) &&= \sigma_x^2 \mathbf{h}\mathbf{h}^T + \sigma_z^2 \mathbf{I} \end{aligned}$$

- Misclassification probability:

$$P_M^{\delta_{\text{Bayes}}} = Q \left(\frac{m_2 - m_1}{2 \sqrt{\sigma_x^2 + \frac{\sigma_z^2}{\|\mathbf{h}\|^2}}} \right)$$



Fixed difficulty level \rightarrow
 $m_1, m_2, \sigma_x, \sigma_z$ & $\|\mathbf{h}\|$ fixed!

2. Unstructured: Full ML

- Recall full Bayes classifier:

$$\delta_{Bayes}(\mathbf{y}) = \mathbf{1}\{\Delta^T \Sigma^{-1}(\mathbf{y} - \mu) \geq 0\}$$

- Unstructured full ML:

- Use $\delta_{Bayes}(\mathbf{y})$ with ML estimates of $\hat{\Delta}$, $\hat{\Sigma}$, $\hat{\mu}$

- Empirical Fisher rule:

$$\delta_{Unstructured}(\mathbf{y}) = \mathbf{1}\{\hat{\Delta}^T \hat{\Sigma}^{-1}(\mathbf{y} - \hat{\mu}) \geq 0\}$$

2. Unstructured: Full ML (continued)

Misclassification probability:

$$P_M^{\delta_{Unstructured}} = Q\left(\frac{\hat{\Delta}^T \hat{\Sigma}^{-1} \hat{\Delta}}{2(\hat{\Delta}^T \hat{\Sigma}^{-1} \Sigma \hat{\Sigma}^{-1} \hat{\Delta})^{1/2}}\right)$$

Theorem 1

If $n/p \rightarrow 0$ and $\|\mathbf{h}\| = \text{constant}$, then $P_M^{\delta_{Unstructured}} \xrightarrow{p} 1/2$

Ignoring structure and estimating all parameters is not a good strategy for inference in high-dimensional scenarios with few samples

3. Full Structure

- Recall sensing model:

$$\mathbf{Y}_i = \mathbf{h}X_i + \mathbf{Z}_i$$

- Structure-based Full Structure:

- **h** KNOWN

- Project data onto **h**, then use ML estimates of parameters

- Projected Empirical Fisher rule:

$$\delta_{Full-Structure}(\mathbf{y}) = \mathbf{1}\left\{\mathbf{h}^T(\mathbf{y} - \hat{\boldsymbol{\mu}})\text{sign}(\hat{\Delta}^T \mathbf{h}) \geq 0\right\}$$

3. Full Structure (continued)

Misclassification probability:

$$P_M^{\delta_{Full-Structure}} = \frac{1}{2} \sum_{i=1}^2 Q \left((-1)^{i-1} \frac{\mathbf{h}^T (\hat{\boldsymbol{\mu}} - \boldsymbol{\mu}_i) \text{sign}(\hat{\Delta}^T \mathbf{h})}{(\mathbf{h}^T \boldsymbol{\Sigma} \mathbf{h})^{1/2}} \right)$$

Theorem 2

If $n \rightarrow \infty$ and $\|\mathbf{h}\| = \text{constant}$, then $P_M^{\delta_{Full-Structure}} \xrightarrow{p} P_M^{\delta_{Bayes}}$

Full knowledge of low-dimensional structure knowledge yields Bayes-optimal performance.

4. Structured ML

- Recall sensing model:

$$\mathbf{Y}_i = \mathbf{h}X_i + \mathbf{Z}_i$$

- Structured ML:

- Sensing model known but **h UNKNOWN**
- Use ML estimate of $\Delta = (m_2 - m_1)\mathbf{h}$ as a proxy for **h**
 - Recall: if **h** known, $P_{\text{error}} \rightarrow P_{\text{Bayes}}$

- Projected Empirical Fisher rule:

$$\delta_{\text{Structured-ML}}(\mathbf{y}) = \mathbf{1}\left\{\hat{\Delta}^T(\mathbf{y} - \hat{\boldsymbol{\mu}}) \geq 0\right\} \quad \leftarrow \text{Essentially, Naive Bayes}$$

4. Structured ML (continued)

Misclassification probability:

$$P_M^{\delta_{Structured-ML}} = \frac{1}{2} \sum_{i=1}^2 Q \left((-1)^{i-1} \frac{\hat{\Delta}^T (\hat{\boldsymbol{\mu}} - \boldsymbol{\mu}_i)}{(\hat{\Delta}^T \boldsymbol{\Sigma} \hat{\Delta})^{1/2}} \right)$$

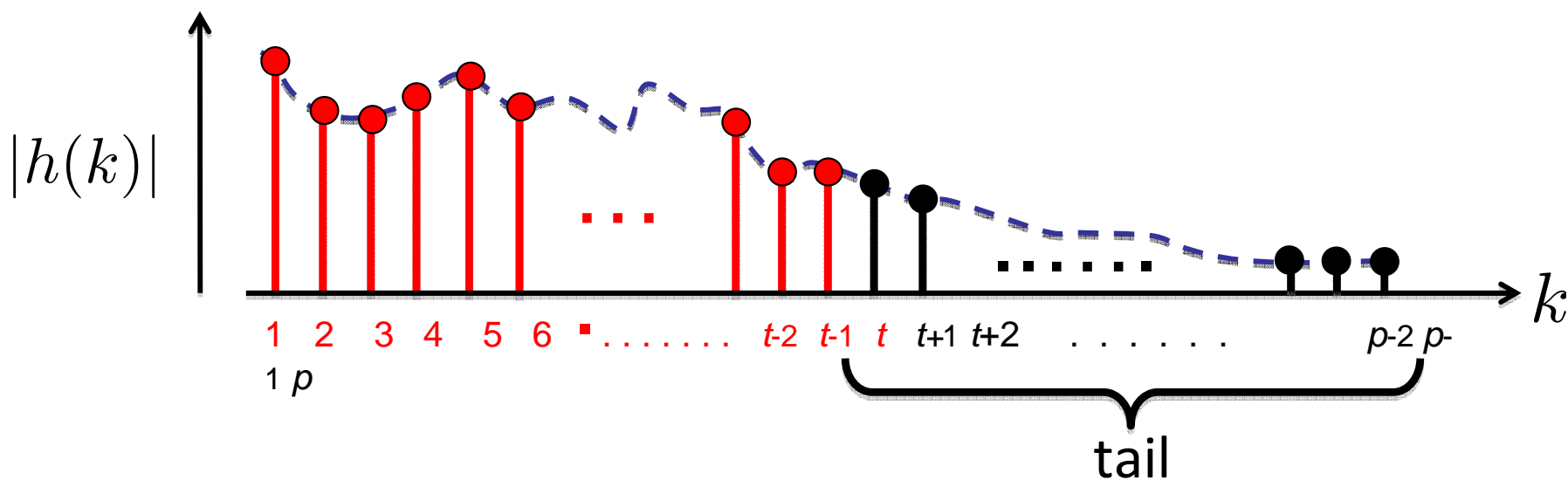
Theorem 3

If $n/p \rightarrow 0$ and $\|\mathbf{h}\| = \text{constant}$, then $P_M^{\delta_{Structured-ML}} \xrightarrow{p} 1/2$

- Simple ML estimation of low-dimensional structure does not work.
- Additional structure is needed!

5. Structured Sparsity

- “Sparsity” assumption: energy in \mathbf{h} localized in few components



- Tail energy goes to zero: As $t, p \rightarrow \infty$,
$$\sum_{k=t+1}^p h^2(k) \rightarrow 0$$

5. Structured Sparsity (continued)

■ Truncation estimator
$$\hat{\Delta}_t(j) = \begin{cases} \hat{\Delta}(j) & \text{if } j \leq t \\ 0 & \text{else} \end{cases}$$

Theorem 4

(a) If $p, n, t \rightarrow \infty$ and $n/t \rightarrow \infty$, then $E[\|\hat{\Delta}_t - \Delta\|^2] \rightarrow 0$

(b) If $n/p \rightarrow 0$ and $\|\mathbf{h}\| = \text{constant}$, then $P_M^{\delta_{\text{Structured-Sparsity}}} \xrightarrow{p} P_M^{\delta_{\text{Bayes}}}$

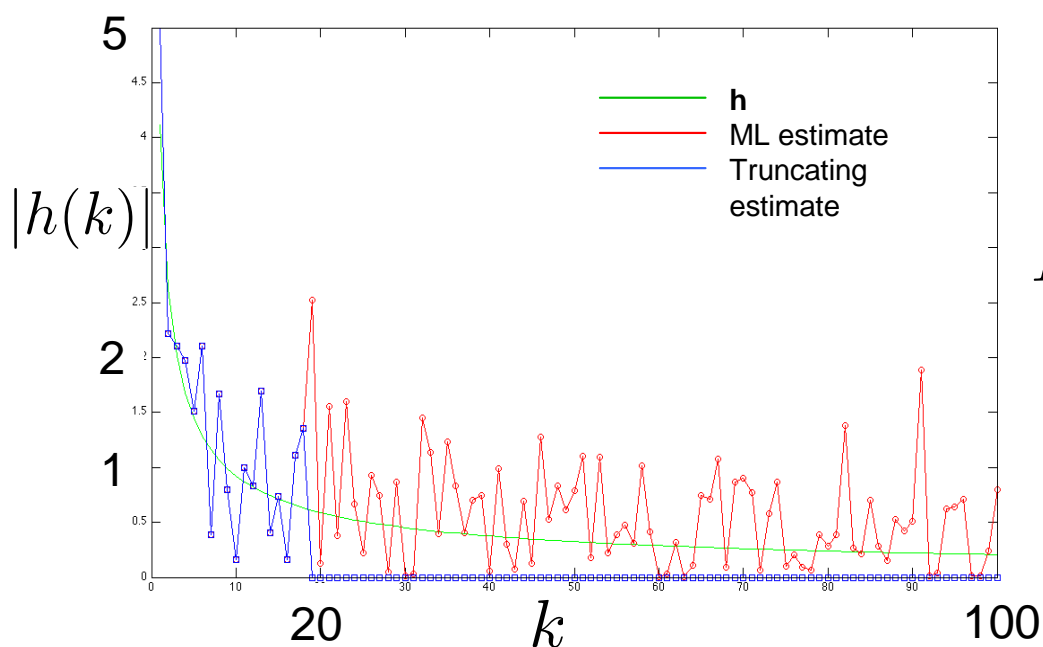
■ Knowledge of sensing structure + sparsity \rightarrow

Asymptotically optimum performance!

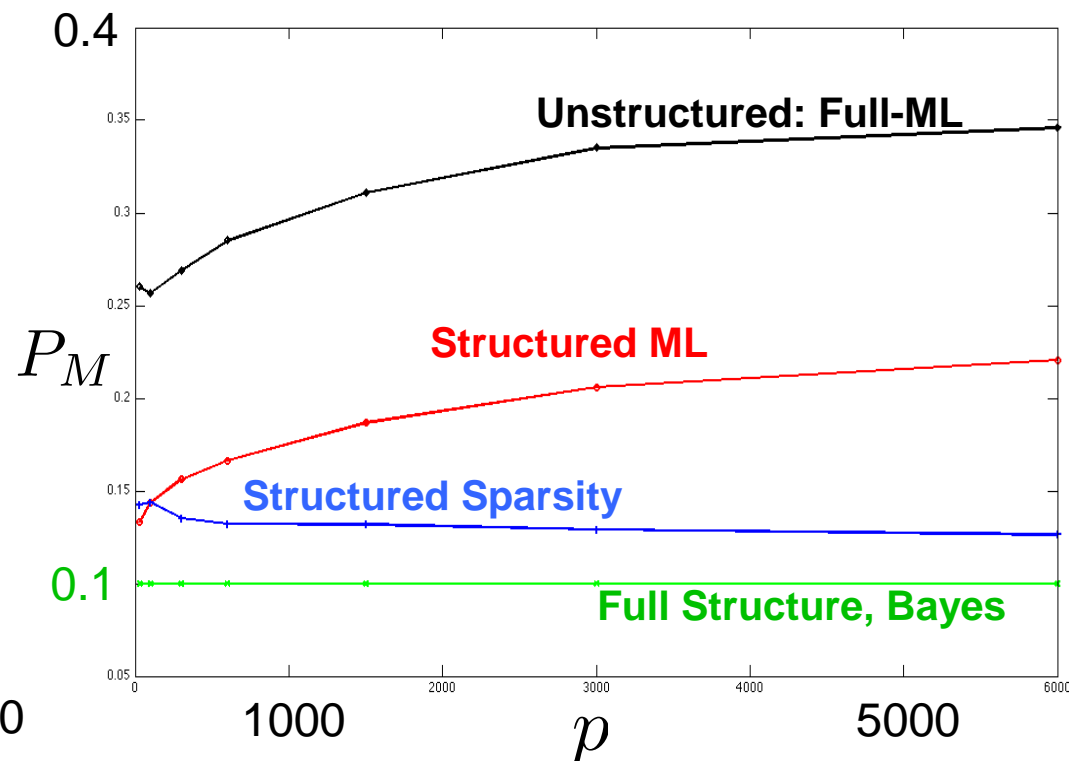
Simulation – Polynomially decaying \mathbf{h}

$$P_M^{\delta_{Bayes}} = 0.1, \quad \|\mathbf{h}\| = 2, \quad \sigma_x = 1, \quad \sigma_z = 2, \quad n = p^{0.7}, \quad t = n^{0.95}$$


Polynomially decaying \mathbf{h} and its
ML & truncating estimates



Misclassification Probability



Concluding remarks

- In many problems of engineering interest
 - *# samples* \ll *data-dimension*
 - there exists a latent low-dimensional sensing structure
- Totally ignoring sensing structure \rightarrow no better than guessing
- General knowledge of sensing structure alone, however, is insufficient
- Knowledge of sensing structure
+
Additional structure in \mathbf{h}  Optimum asymptotic
classification
performance