

Problems and assumptions

- Outliers are omnipresent : economic science, astronomical science,... computer vision, etc..
- Statistics of the hat matrix used for outlier identification. We propose a nonlinear extension of the standard hat matrix: The **kernel hat matrix**
- A model selection criterion is derived for **dominant data subset extraction**
- Performances of the proposed approaches are studied on simulated and real data sets.

From the hat matrix to the kernel hat matrix

1 Hat matrix

Let $Z = [X:y] = [z_1, z_2, \dots, z_n]^T$ is the $n \times (d+1)$ matrix that combines both the matrix of predictors X and the response vectors y . The augmented hat matrix is: $\mathcal{H} = Z(Z^T Z)^{-1} Z^T$

2 « Kernelized » Hat matrix

Consider a nonlinear transformation $\phi: \mathbb{R}^{d+1} \rightarrow \mathbb{R}^p$ and the $n \times p$ matrix $\Phi = (\phi(z_1), \phi(z_2), \dots, \phi(z_n))^T$

$$\mathcal{H}_\Phi = \Phi \underbrace{(\Phi^T \Phi)^{-1} \Phi^T}_{p \times p \text{ (high dimension)}} \xrightarrow{\text{Regularization parameter}} \mathcal{H}_\Phi = \Phi \Phi^T \underbrace{(\Phi \Phi^T + \gamma I)^{-1}}_{n \times n \text{ (low dimension)}} \xrightarrow{\mathcal{K} = \Phi \Phi^T}$$

$$\mathcal{H}_K = \mathcal{V} \Omega \mathcal{V}^T \xleftarrow{\text{Eigenvalue decomposition: } \mathcal{V} = (V_1, V_2, \dots, V_n) \text{ matrix of eigenvectors of } \mathcal{K}, \Lambda = (\lambda_1, \lambda_2, \dots, \lambda_n) \text{ matrix of eigenvalues of } \mathcal{K}} \mathcal{H}_K = \mathcal{K} (\mathcal{K} + \gamma I)^{-1}$$

Effective degree of freedom $\Omega = \text{diag}\left(\frac{\lambda_1}{\lambda_1 + \gamma}, \frac{\lambda_2}{\lambda_2 + \gamma}, \dots, \frac{\lambda_n}{\lambda_n + \gamma}\right)$

4 Proposed model selection for dominant data set extraction:

- Goal:** optimizing the separation between the h_{ii} 's distribution of the dominant data subset from the one of the outliers.

- Decision threshold:** $\mathcal{H}_{Kii} \leq \text{trace}(\mathcal{H}_K)/n \quad 1 < i < n$

This threshold gives an initial partition of the h_{ii} 's distribution into two distributions :

- 1) The « Dominant » distribution $\mathcal{H}^{(D)}$ of mean $\mu^{(D)}$ and variance $\sigma^{(D)}$.
- 2) The « Outlier » distribution $\mathcal{H}^{(O)}$ of mean $\mu^{(O)}$ and variance $\sigma^{(O)}$.

- Proposed model selection:** We consider the couple (σ^*, γ^*) which maximizes the ratio between-subset variance (i.e separation) and within-subset variance (i.e overlap):

$$\delta(\sigma, \gamma) = \frac{(\mu^{(D)}(\sigma, \gamma) - \mu^{(O)}(\sigma, \gamma))^2}{\sigma^{(D)}(\sigma, \gamma) + \sigma^{(O)}(\sigma, \gamma)} \quad \text{Linear Fisher's discriminant}$$

- How does \mathcal{H}_K evolve with respect to (σ, γ) ???** We show that a pair (σ, γ) exists which best separates a dominant subpopulation from outliers !!! Look at this example of three unidimensional data points

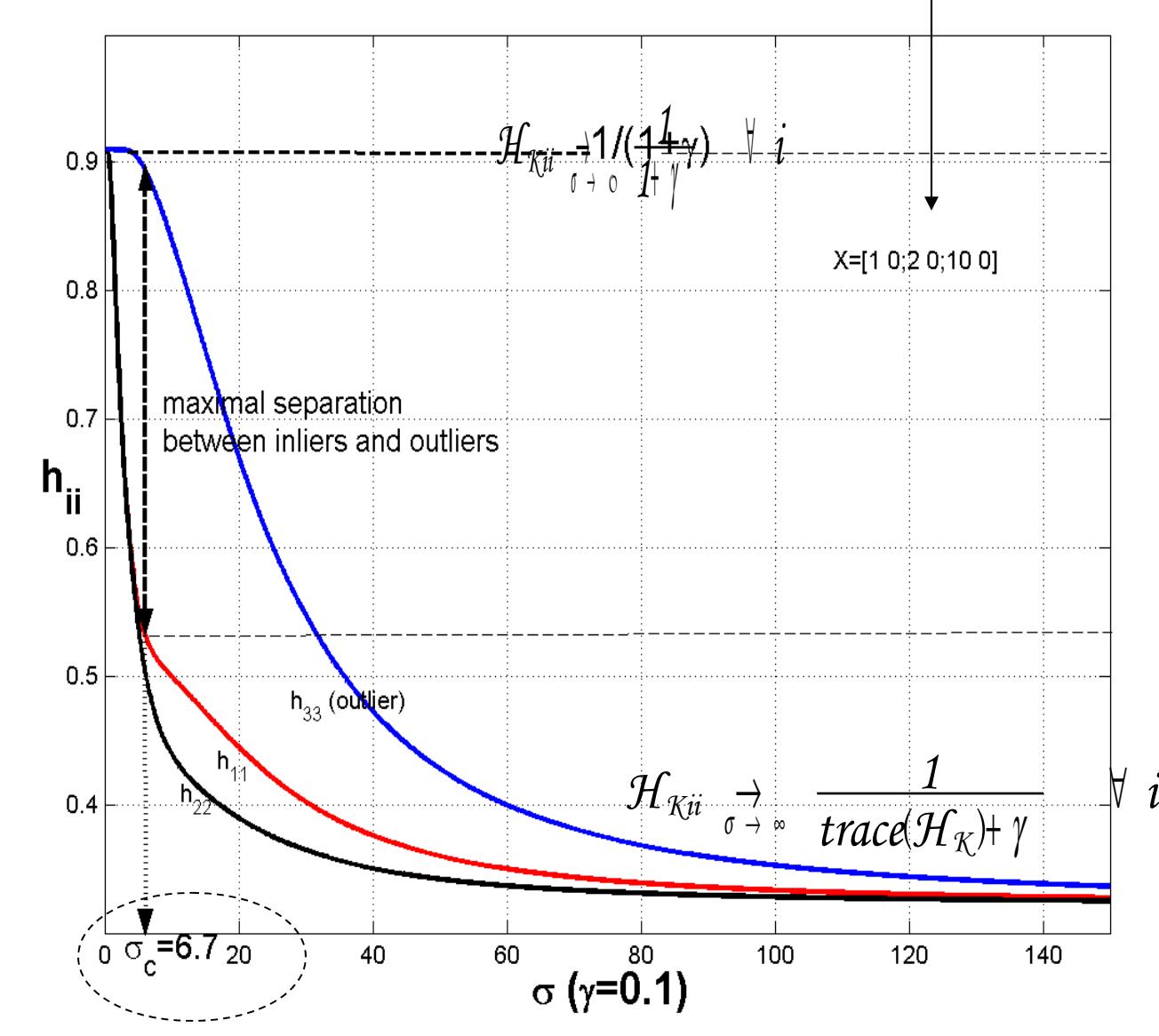
A dominant subpopulation of two points ($Z1$ and $Z2$) and one outlier ($Z3$)

The evolution of the corresponding h_{ii} with respect to σ shows a maximal separation between h_{33} and the pair (h_{11}, h_{22}) for a critical value of σ (σ_c in the picture on the right) !!!

An analytical expression of σ_c can be obtained

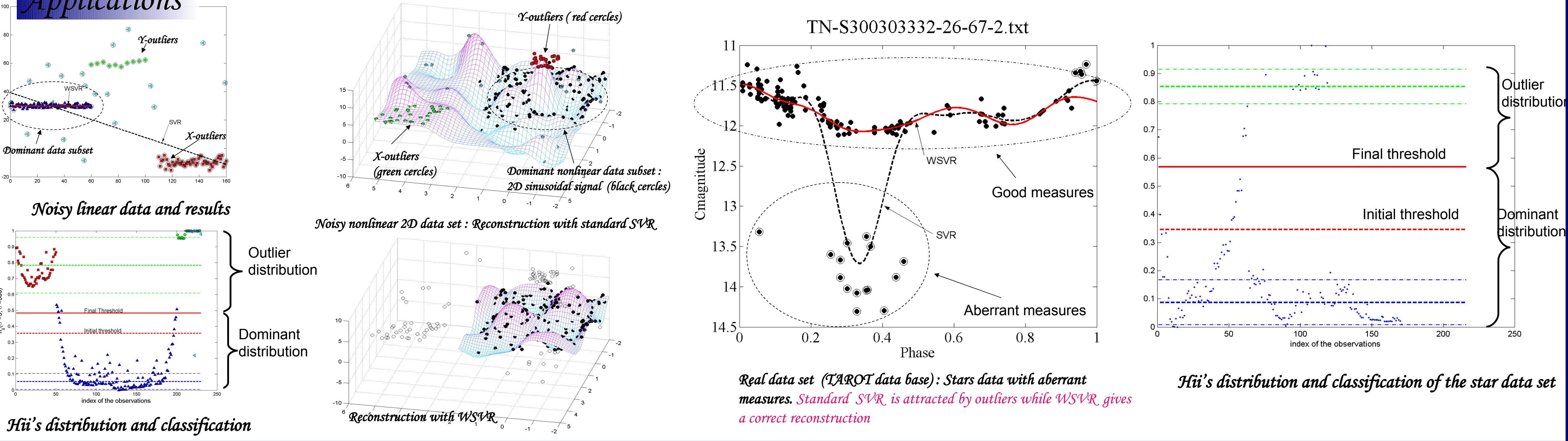
$$\sigma_c = \frac{\|Z_1 - Z_2\|^2}{2} \frac{\alpha^2 - 1}{\ln \alpha}$$

$$\text{with } \alpha = \frac{\|Z_i - Z_c\|^2}{\|Z_1 - Z_2\|^2}$$



5 Conclusion : The proposed model selection is used as learning stage for data classification before regression task: application to support vector regression : the proposed approach is called: Weighted SVR (WSVR)

Applications



Conclusion

- A new robust learning strategy for outlier detection has been proposed.
- The Gaussian Kernel hat matrix presents discriminative properties under the condition to choose appropriate values for kernel parameters.

Futur work

- Automatic selection of kernel parameters
- Spectral analysis of the kernel hat matrix for multiple modality extraction
- Application to computer vision problems : optical flow estimation, segmentation