

NASCA '18

2-6 July 2018

Kalamata, Greece

CONSTRAINED WEIGHTED FEATURE SELECTION

*NASCA 2018: Numerical Analysis and Scientific Computation with
Application Conference, 2018*

Denis Hamad

Samah Hijazi

Mariam Kalakech

Ali Kalakech

ulco UNIVERSITÉ
DU LITTORAL
CÔTE D'OPALE



Région
Hauts-de-France

LISIC
Laboratoire d'informatique
Signal & Image de la Côte d'Opale



GROUPEMENT DE RECHERCHE
EN AUTOMATISME, INTELIGENCE
ET SYSTEMES HUMAINS-MACHINES



المجلس الوطني للبحوث العلمية

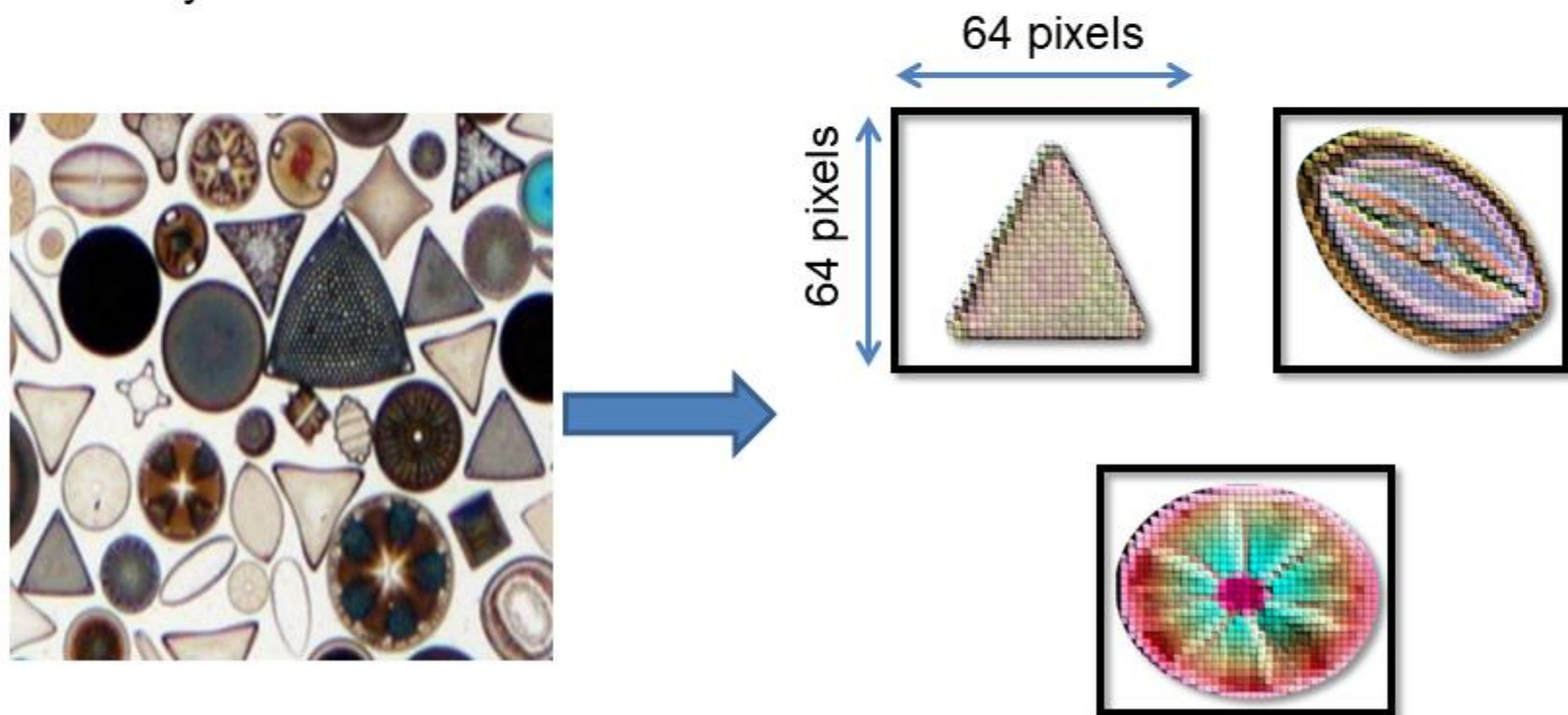
National Council for Scientific Research



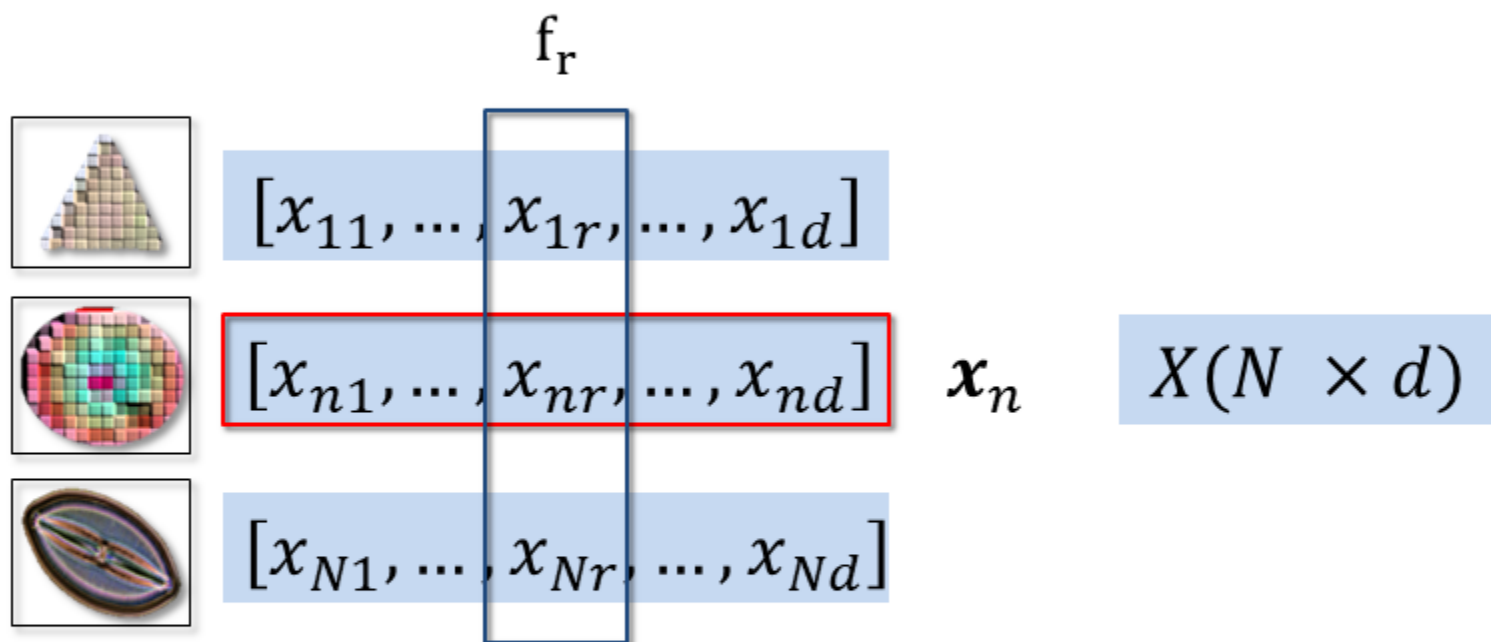
Université Libanaise
Ecole Doctorale
Sciences et Technologie

Context, Motivation & Problem

- In many machine learning applications the data used is provided with a very large number of features.
- Where only few are relevant and not redundant.

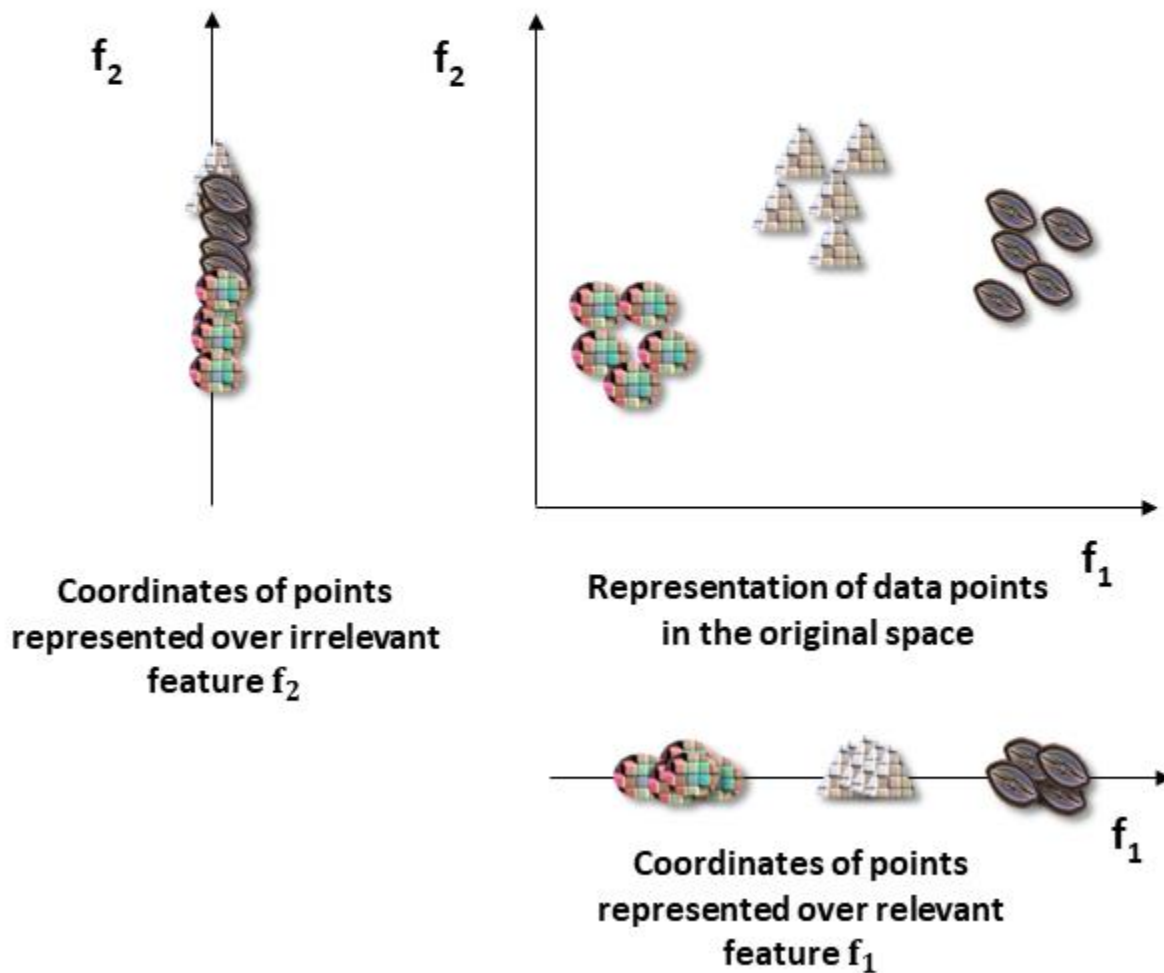


Context, Motivation & Problem



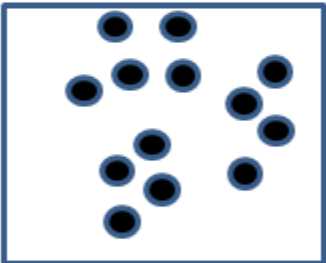

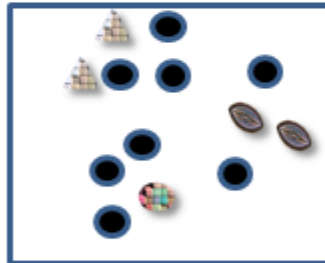
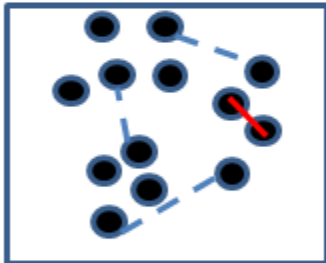
- Each object \mathbf{x}_n ($n= 1 \dots, N$) is characterized by a large number of features f_r ($r =, 1 \dots, d$).
- x_{nr} represents the value of the n -th data object on the r -th feature.
- The performance of machine learning algorithms might be degraded when applied on such high dimensional data.

Solution: Feature Selection



Feature Selection:

Learning Contexts

Unsupervised Learning	Supervised Learning	Semi-supervised Learning	Constrained Learning
 <p data-bbox="142 991 494 1165">We know nothing about the data labels</p> <ul data-bbox="142 1236 504 1315" style="list-style-type: none"> • <i>Laplacian Score</i> • <i>Variance Score</i> 	 <p data-bbox="562 991 938 1165">We know the data labels very well</p> <ul data-bbox="581 1236 942 1315" style="list-style-type: none"> • <i>Fisher Score</i> • <i>Relief Algorithm</i> 	 <p data-bbox="1051 958 1306 1196">We know something about data labels</p> <ul data-bbox="987 1219 1309 1362" style="list-style-type: none"> • <i>Semi-supervised Laplacian Score</i> • <i>Semi-supervised logistic I-Relief</i> 	 <p data-bbox="1441 958 1779 1196">We have pairwise constraints on data</p> <ul data-bbox="1479 1236 1725 1268" style="list-style-type: none"> • <i>Simba-Sc</i> <div data-bbox="1464 1296 1773 1368" style="border: 1px solid black; padding: 2px; display: inline-block;"> Relief-Sc </div>

Outline

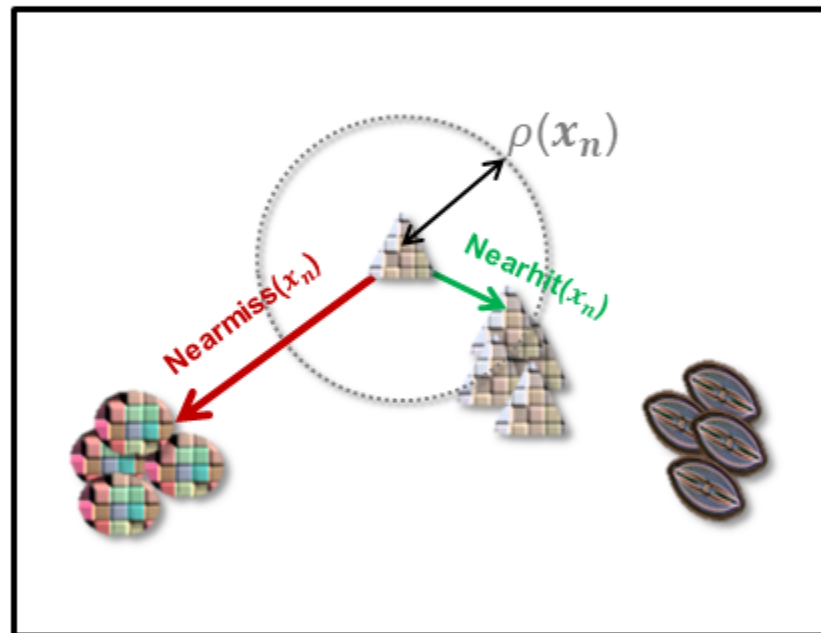
- Hypothesis-Margin: Relief-Sc (Relief with Side Constraints)
- Algorithmic Comparison
- Experimental results on Relief-Sc
- Constraint Selection
- Constraint propagation
- Conclusion & Future Work

Hypothesis-Margin

Supervised Context

Kira K, Rendell LA (1992)

- **Nearmiss** of an instance x_n is the nearest sample to x_n with a different label
 - **Nearhit** of an instance x_n is the nearest sample to x_n with the same label
- **Hypothesis Margin** of x_n denoted $\rho(x_n)$ is the **difference** between the distance to its **nearmiss** and the distance to its **nearhit**.
 - **Hypothesis Margin** is the largest distance an instance of the dataset can travel without altering the labeling of instances.



Gilad-Bachrach R., Navot A., Tishby (N. 2004)

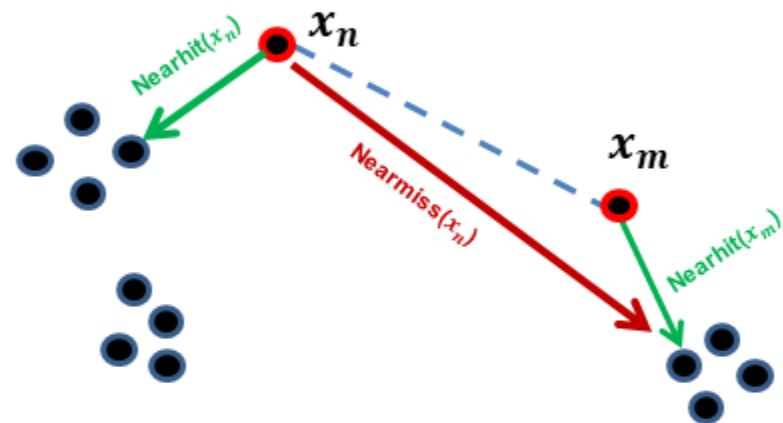
$$\rho(x_n) = \text{diff}(x_n, \text{Nearmiss}(x_n)) - \text{diff}(x_n, \text{Nearhit}(x_n))$$

Hypothesis-Margin

Constrained Context

Yang M., Song J. (2010)

- **Nearmiss** of an instance x_n is the nearest sample to x_m
- **Nearhit** of an instance x_n is the nearest sample to x_n



$$\rho(x_n, x_m) = \text{diff}(x_n, \text{Nearhit}(x_m)) - \text{diff}(x_n, \text{Nearhit}(x_n))$$

Constrained Weighted Feature Selection

- Weighted Hypothesis Margin of $(\mathbf{x}_n, \mathbf{x}_m)$

$$\rho(\boldsymbol{\omega}, (\mathbf{x}_n, \mathbf{x}_m)) = \boldsymbol{\omega}^T [\text{diff}(\mathbf{x}_n, \text{Nearhit}(\mathbf{x}_m)) - \text{diff}(\mathbf{x}_n, \text{Nearhit}(\mathbf{x}_m))]$$

- Overall Weighted margin

$$\rho = \boldsymbol{\omega}^T \sum_{(\mathbf{x}_n, \mathbf{x}_m) \in \mathcal{C}} [\text{diff}(\mathbf{x}_n, \text{Nearhit}(\mathbf{x}_m)) - \text{diff}(\mathbf{x}_n, \text{Nearhit}(\mathbf{x}_m))]$$

Sun Y. and Li J. (2006)

Z

- Problem Formulation:

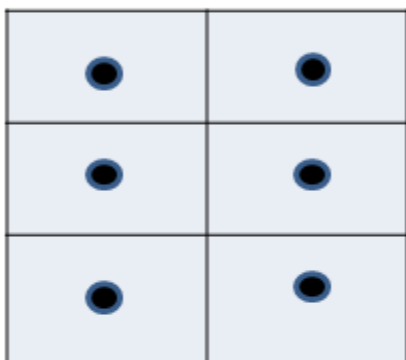
- We want to find the weight vector $\boldsymbol{\omega}$ that maximizes the overall margin
- Then the Problem is defined as

$$\max_{\boldsymbol{\omega}} \boldsymbol{\omega}^T \mathbf{Z} \quad \text{s.t. } \|\boldsymbol{\omega}\|^2 = 1 \text{ and } \boldsymbol{\omega} \geq 0$$

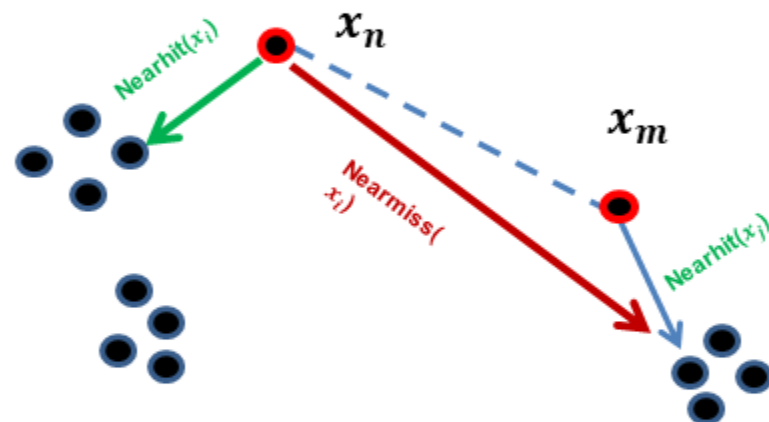
$$\text{where } \boldsymbol{\omega} = (\omega_1, \dots, \omega_r, \dots, \omega_d)$$

Relief-Sc

(Relief with Side Constraints)



Cannot-link Constraints Set



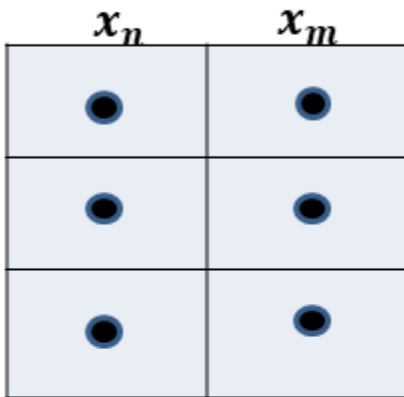
$$\begin{array}{c}
 \mathbf{x}_n \\
 \mathbf{x}_m \\
 \vdots \\
 \mathbf{x}_N
 \end{array}
 = \begin{bmatrix}
 \mathbf{f}_1 & \cdots & \mathbf{f}_r & \cdots & \mathbf{f}_d \\
 x_{11} & \cdots & x_{1r} & \cdots & x_{1d} \\
 \cdots & \cdots & \cdots & \cdots & \cdots \\
 x_{n1} & \cdots & x_{nr} & \cdots & x_{nd} \\
 \cdots & \cdots & \cdots & \cdots & \cdots \\
 x_{N1} & \cdots & x_{Nr} & \cdots & x_{Nd}
 \end{bmatrix}$$

$$\mathbf{z} = (0, \quad 0 \quad \cdots \quad 0)$$

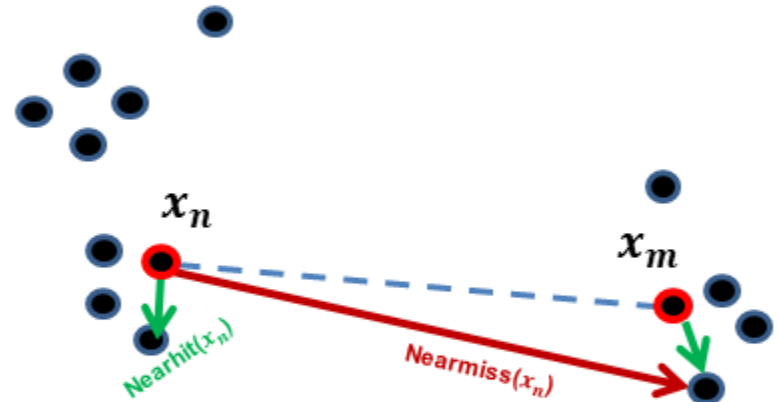
$$z_r = z_r + |x_{nr} - \text{Nearhit}(x_{mr})| - |x_{nr} - \text{Nearhit}(x_{nr})|$$

Relief-Sc

(Relief with Side Constraints)



Cannot-link Constraints Set



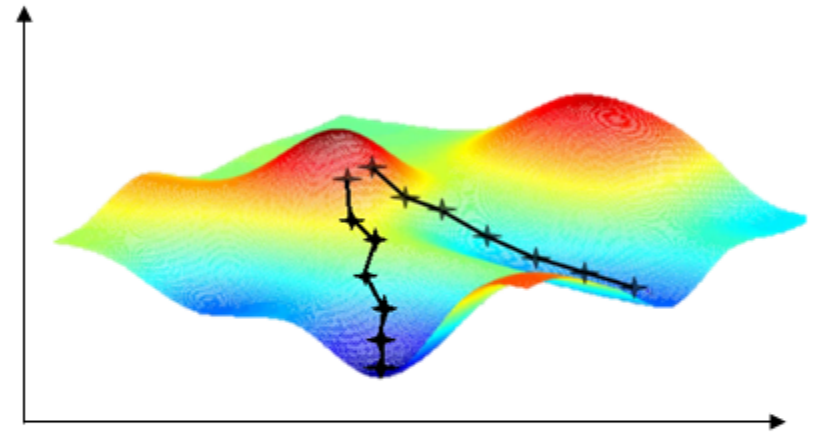
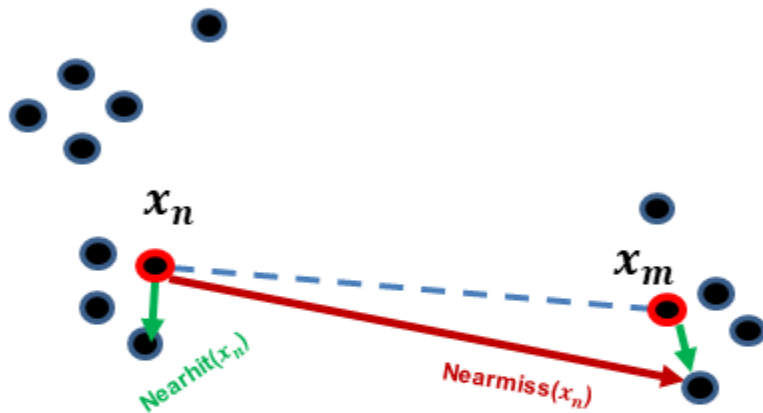
$$z_r = z_r + |x_{nr} - \text{Nearhit}(x_{mr})| - |x_{nr} - \text{Nearhit}(x_{nr})|$$

$$\mathbf{z}^+ = [\max(z_1, 0), \dots, \max(z_d, 0)]^T$$

$$\boldsymbol{\omega} = \frac{\mathbf{z}^+}{\|\mathbf{z}^+\|}$$

\mathbf{f}_1	\mathbf{f}_r	\mathbf{f}_d
$[\omega_1$	ω_r	$\omega_d]$

Algorithmic comparison: Simba-Sc



Yang M., Song J. (2010)

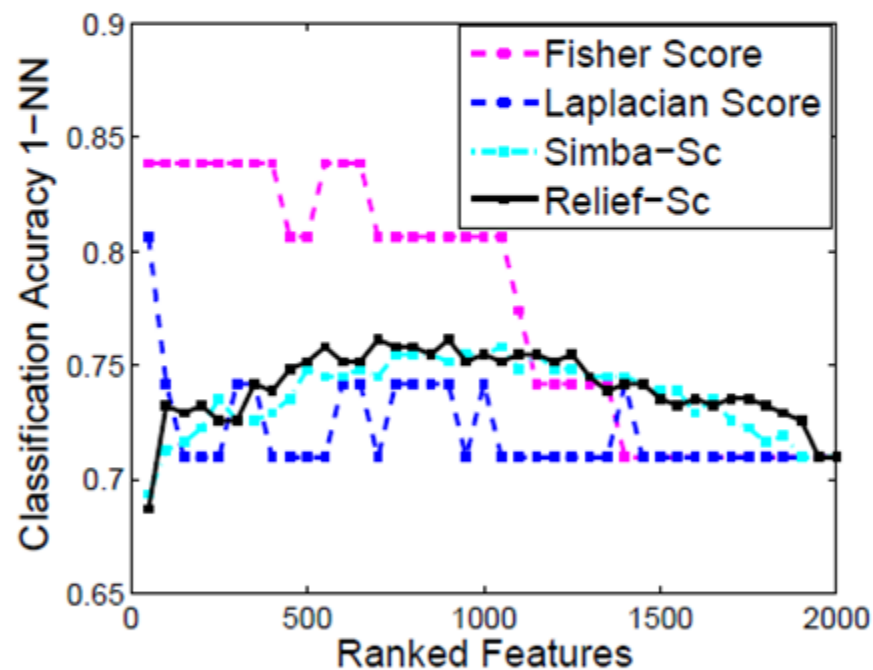
$$\omega_r = \omega_r + \frac{1}{2} \left[\frac{(x_{nr} - \text{Nearhit}(x_{mr}))^2}{\|x_n - \text{Nearhit}(x_m)\|_\omega} - \frac{(x_{nr} - \text{Nearhit}(x_{nr}))^2}{\|x_n - \text{Nearhit}(x_n)\|_\omega} \right] \omega_r$$

$$\omega = \frac{\omega}{\|\omega\|_\infty}$$

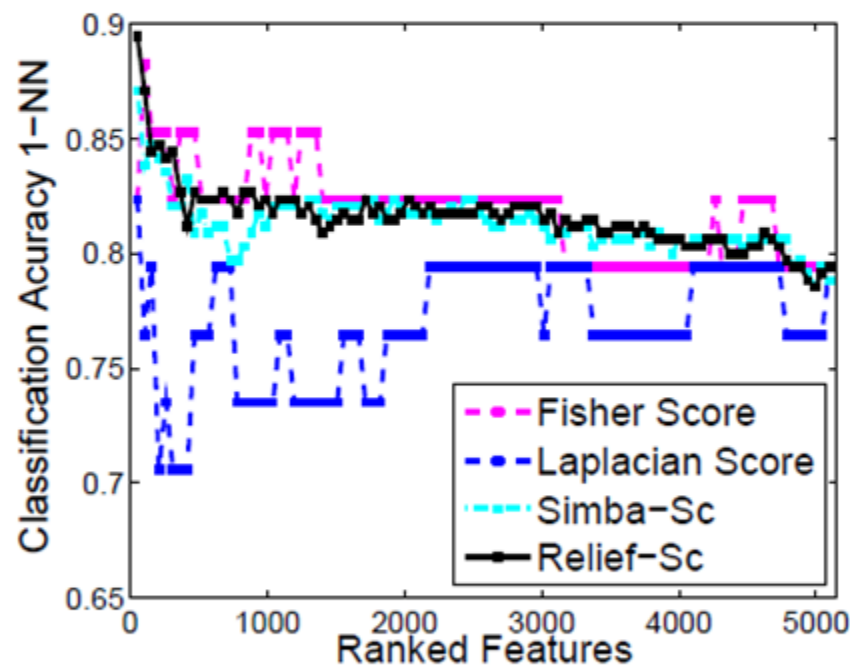
f_1	f_r	f_d
$[\omega_1$	ω_r	$\omega_d]$

Experimental Results:

Classification Accuracy of Different Feature Selection Algorithms



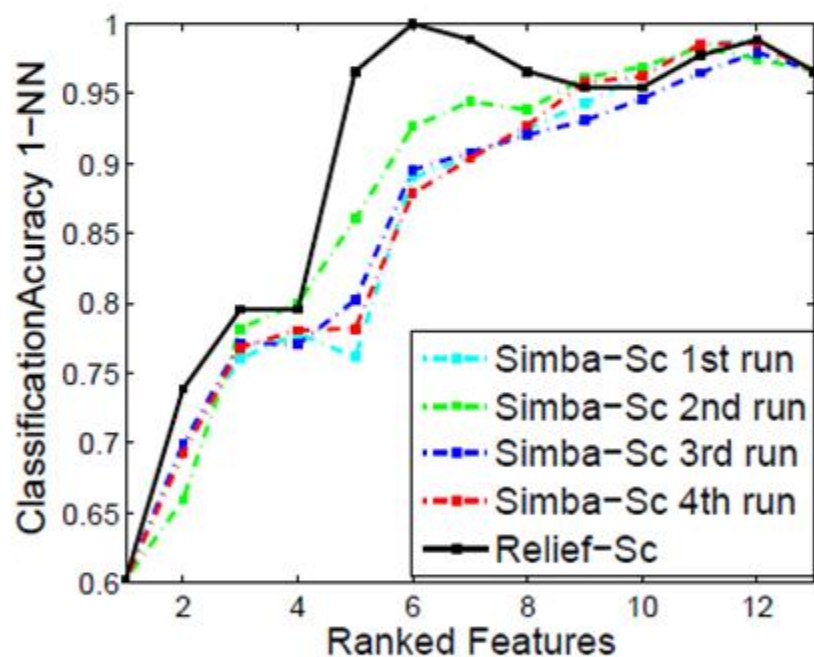
(a) ColonCancer



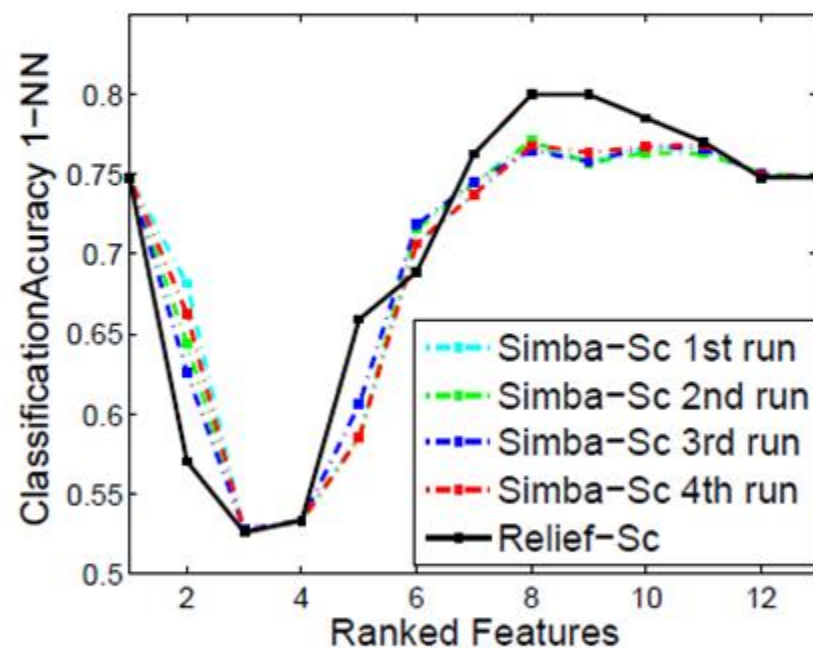
(b) Leukemia

Experimental Results:

Solutions Comparison of Relief-Sc and Simba-Sc



(a) Wine



(b) Heart

Selection of Cannot-link Constraint Set

x_n	x_m
●	●
●	●
●	●

Cannot-link Constraints Set



1. Find the Laplacian matrix L of the data
2. Find eigenvectors V and eigenvalues λ of $L=D-S$

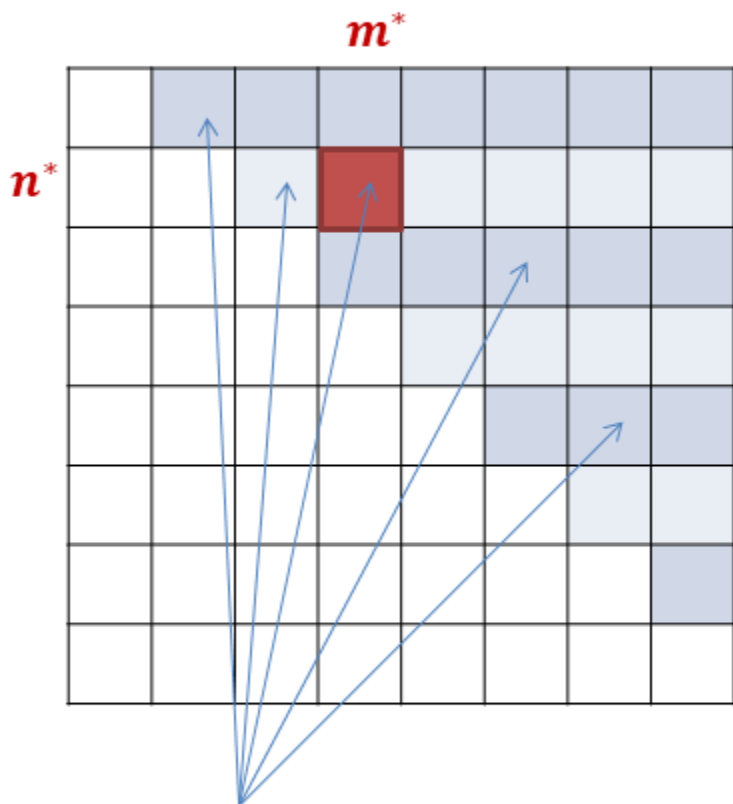
3. Find the point of minimum value on the second eigenvector v_2



4. Find the data couple that can most change the position of the minimum value $v_2(i^*)$ on v_2 such that:

$$(x_{n^*}, x_{m^*}) = \operatorname{argmax}_{n, m \in \{1 \dots N\}} \left| \frac{dv_2(i^*)}{ds_{nm}} \right|$$

Selection of Cannot-link Constraint Set



Sensitivities of (x_n, x_m)

We find the sensitivity of each couple and then The maximum (x_{n^*}, x_{m^*}) among all.

4. Thus, we calculate the sensitivity of $v_2(i^*)$ to each couple in the dataset and store them in a matrix.

Now, we find the data having the highest sensitivity

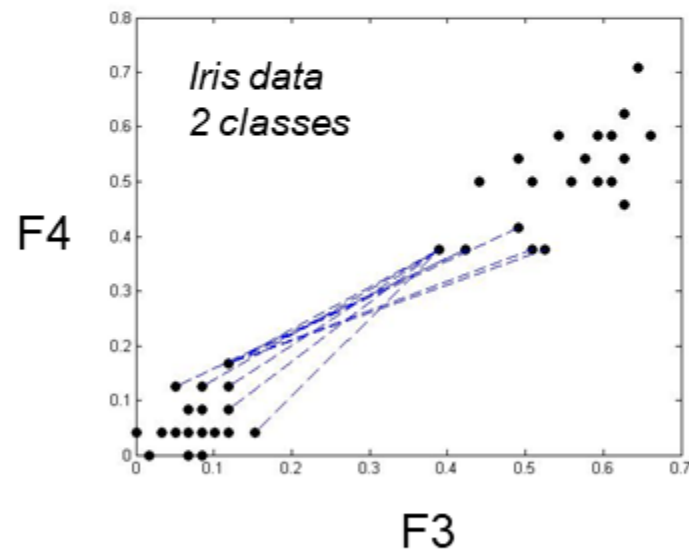
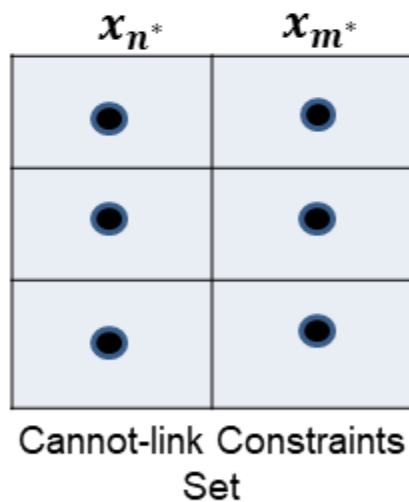
$$(x_{n^*}, x_{m^*}) = \operatorname{argmax}_{n, m \in \{1 \dots N\}} \left| \frac{dv_2(i^*)}{ds_{nm}} \right|$$

Where

$$\begin{aligned} \left| \frac{dv_2(i^*)}{ds_{nm}} \right| &= \left| \sum_{p>2}^N \frac{v_2^T [\partial \mathbf{L} / \partial s_{nm}] v_p}{\lambda_2 - \lambda_p} v_{p(i^*)} \right| \\ &= \left| \sum_{p>2}^N \frac{v_2^T [(e_n - e_m)(e_n - e_m)^T] v_p}{\lambda_2 - \lambda_p} v_{p(i^*)} \right| \end{aligned}$$

Selection of Cannot-link Constraint Set

5. After finding the indexes (n^*, m^*) , We Actively query for a constraint on this particular couple.
6. Thus, obtaining the constraints of highest utility.

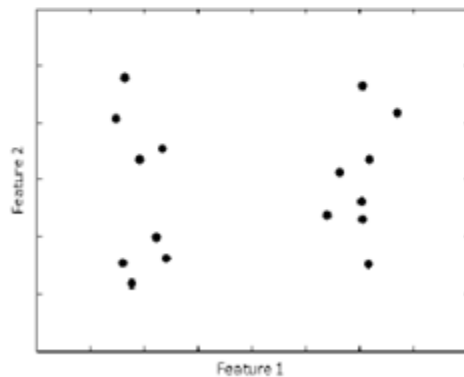


However, we can't ask the user a lot of questions, so the number of constraints we can obtain is limited

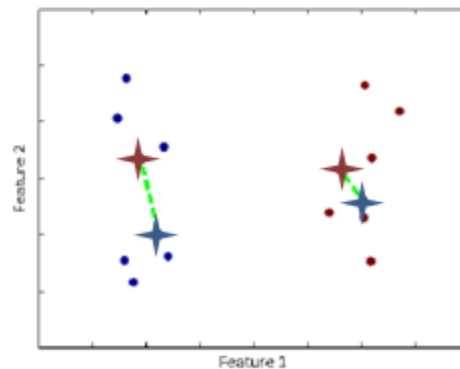
Constraints Propagation

- Relief family margin-based algorithms need a set of points that is sufficiently big to calculate the margin.
- we query an oracle for constraints
- we might not be able to ask for a large number of queries

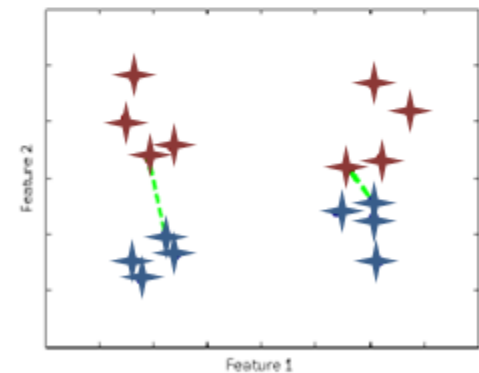
Therefore we needed a constraint propagation method



(a) Data



(b) Instance-level constraints



(c) Propagated information

Constraints Propagation

- We initialize the constraints matrix Q_{nm} (from previous constraints selection step) as follows:

$$Q_{nm} = \begin{cases} 1, & \text{if } (x_n, x_m) \in \mathcal{C} \\ 0, & \text{otherwise} \end{cases}$$

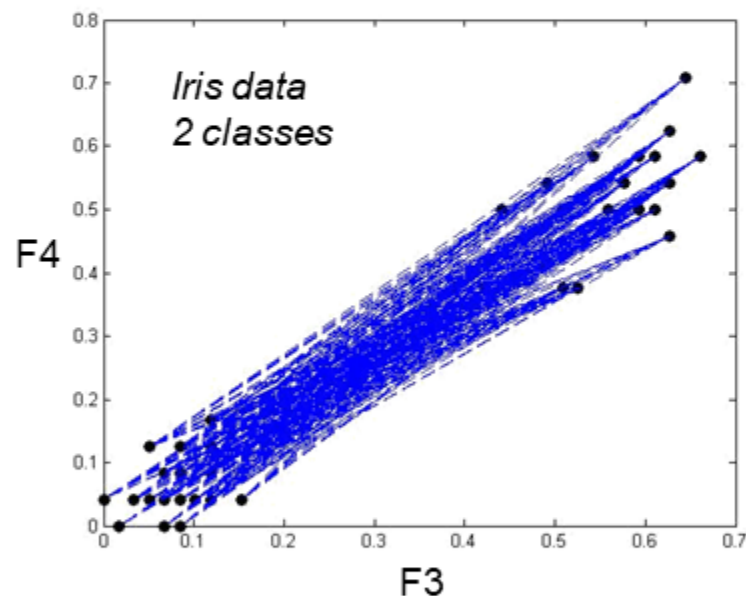
- We calculate the normalized neighborhood similarity graph P_{nm}

$$P_{nm} = \begin{cases} e^{-\frac{\|x_n - x_m\|^2}{2\sigma^2}}, & \text{if } x_m \in k\text{-NN}(x_n) \\ 0, & \text{otherwise} \end{cases}$$

- Finally we use the following matrix completion equation to propagate our constraint on P_{nm} obtaining G^*

$$G^* = (1 - a)^2 (I - aP)^{-1} Q (I - aP)^{-1}$$

Where I is the identity matrix and a is a regularization parameter



Conclusion

- We proposed Relief-Sc, a **weighted feature selection** algorithm that works in a **constrained** environment.
- It is said to find a **unique** relevant feature subset in a **closed-form**.
- We are currently working on the experimental results of Constraints selection and Propagation, we expect to obtain better feature selection with a minimum number of queried constraints.

Acknowledgments

- This work is supported by a scholarship granted by the University of the Littoral Opal Coast (ULCO) and (AUF) in France together with the National Council For Scientific Research in Lebanon (CNRS-L) as a part of ARCUS E2D2 project.

References

- ❑ Gilad-Bachrach R., Navot A., Tishby N. (2004), "Margin based feature selection—theory and algorithms », 21st Int. Conf. on Machine Learning, Canada, pp. 43–50.
- ❑ Kira K, Rendell L.A. (1992) A practical approach to feature selection. In: Proceedings of the ninth int. workshop on Machine learning, pp 249-256.
- ❑ Lu Z., Ip H.H. (2010) Constrained spectral clustering via exhaustive and efficient constraint propagation. In: European Conf. on Computer Vision, Springer, pp 1-14.
- ❑ Sun Y. and Li J. (2006) “Iterative RELIEF for Feature Weighting,” Proc. 23rd Int’l Conf. Machine Learning, pp. 913-920.
- ❑ Wauthier F.L., Jojic N., Jordan M.I. (2012) Active spectral clustering via iterative uncertainty reduction. In: Proceedings of the 18th ACM SIGKDD int. conf. on Knowledge discovery and data mining, ACM, pp 1339-1347.
- ❑ Yang M., Song J. (2010) A novel hypothesis-margin based approach for feature selection with side pairwise constraints. Neurocomputing, 73(16):2859-2872.