



Feature Selection for Classification

Data preprocessing and Feature Selection

MOSAIC PINV15-257 

Miguel García Torres
mgarcia@upo.es

Lenguajes y Sistemas Informáticos, Universidad Pablo de Olavide 

Intelligent Data Analysis research group 

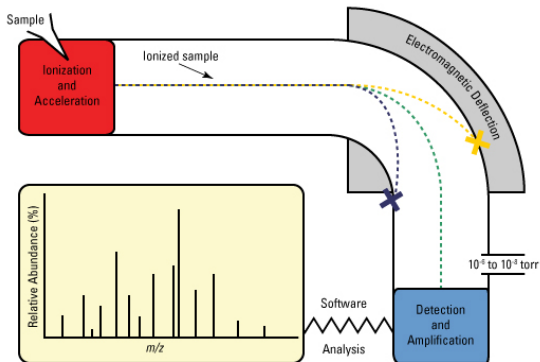
Outline

- 1 **Background**
 - Introduction
 - Feature relevance
 - Feature redundancy
 - Relevance and optimality
- 2 **Feature selection Steps**
 - Feature subset generation
 - Evaluation
- 3 **Feature selection algorithm**
 - Greedy Sequential Search
 - Fast Correlation Based Filter (FCBF)
 - Scatter Search (SS)



Proteomic mass spectrometry analysis

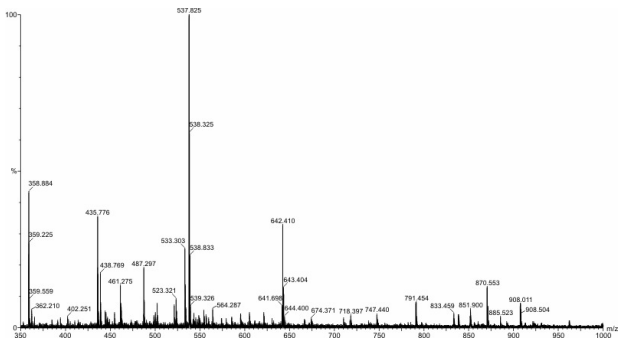
- Enhanced data acquisition \Rightarrow large high-dimensional data!
- MS analysis \Rightarrow Peptide and protein analysis.
- Objective: Biomarker detection.



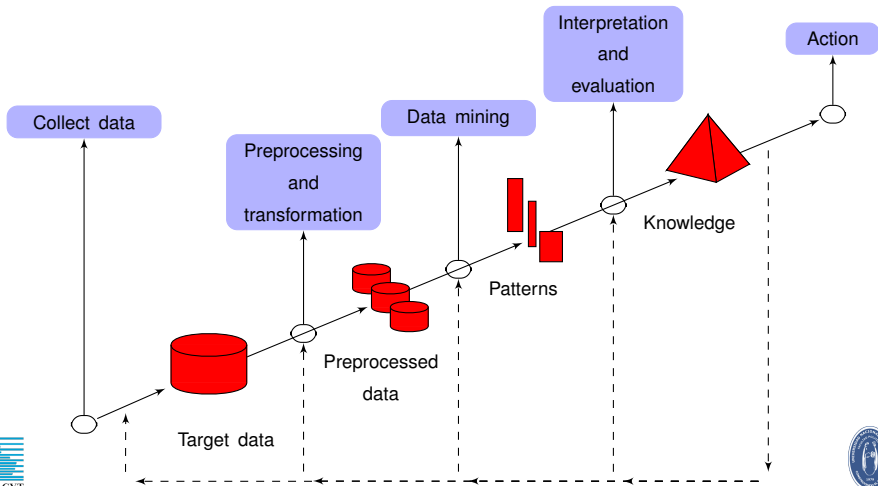
Proteomic mass spectrometry analysis

Challenges

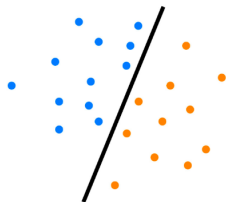
- Small high-dimensional dataset.
- Original signal decomposition unknown.
- No standard data preprocessing workflow.



The Knowledge Discovery in Databases process



Classification



Notation

- $X = \{X_j, j = 1, \dots, d\}$ full set of features.
- $Y \equiv$ the class variable (target class to be learned).
- $E = (\mathbf{x}, y) \equiv$ training set.
- $T = (\mathbf{x}, ?) \equiv$ test set.

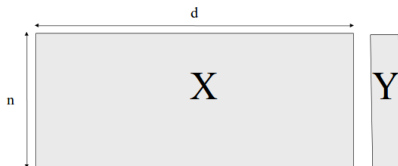


Given E , the aim of classification is to learn a function

$$C : X \rightarrow Y.$$



Feature selection for classification



- Not all the features are equally useful \Rightarrow removing some of them may improve the predictive model \mathcal{C} .
- \Rightarrow The objective of feature selection is to find the subset of features $S \in X$ with which \mathcal{C} achieves the lowest error rate.

Benefits

- Reduction in the cost of acquisition of the data.
- Improvement of the comprehensibility of the model.
- Faster induction of the final classification model.
- Improvement in classification accuracy.

Feature selection (FS)

- FS traditionally focused on finding a highly discriminating power set of features for minimizing the classification error rate.
- Several works have made an effort for defining the different feature types according to their contribution to the meaning of the class concept.
- In this context, feature relevance has arisen as a measure of the amount of relevant information that a feature may contain about the class in classification tasks.
- A feature is considered irrelevant if it contains no information about the class and therefore it is not necessary at all for the predictive task
- Relevant features are those that embody information about the class concept

Feature selection (FS)

- Optimal feature subset defined with respect to the induction algorithm:

*Given an inducer \mathcal{I} , and a training dataset E with features X_1, \dots, X_d , from a distribution \mathcal{D} over the labeled instance space, an **optimal feature subset**, S_{opt} , is a subset of the features such that the accuracy of the induced classifier $\mathcal{C} : \mathcal{I}(\mathcal{D})$ is maximal.*

- Optimal feature subset not necessarily unique.
- Problem: distribution of data unknown.

⇒ Accuracy of the classifier must be estimated from data.

Kohavi & John [19]

$$S_j = X \setminus \{X_j\}.$$

- **Strong relevance** \equiv A feature X_j is strongly relevant iff

$$P(Y|X_j, S_j) \neq P(Y|S_j).$$

- **Weak relevance** \equiv A feature X_j is weakly relevant iff

$$P(Y|X_j, S_j) = P(Y|S_j).$$

and $\exists S'_j$, such that

$$P(Y|X_j, S'_j) \neq P(Y|S_j).$$

- **Irrelevance** \equiv A feature X_j is irrelevant iff

$$\forall S'_j \subseteq S_j, P(Y|X_j, S'_j) = P(Y|S_j).$$

Target concept

$$Y = X_1 \oplus X_2.$$

where

$$X_4 = \bar{X}_2$$

$$X_5 = \bar{X}_3$$

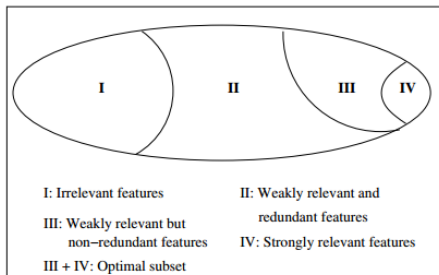
X_1	X_2	X_3	X_4	X_5	Y
0	1	1	0	0	0
0	1	0	0	1	0
0	0	1	1	0	1
0	0	0	1	1	1
1	1	1	0	0	1
1	1	0	0	1	1
1	0	1	1	0	0
1	0	0	1	1	0

- $X_1 \equiv$ strongly relevant.
- $X_2, X_4 \equiv$ weakly relevant.
- $X_3, X_5 \equiv$ irrelevant.
- model with highest accuracy $\{X_1, X_2\}, \{X_1, X_4\}$

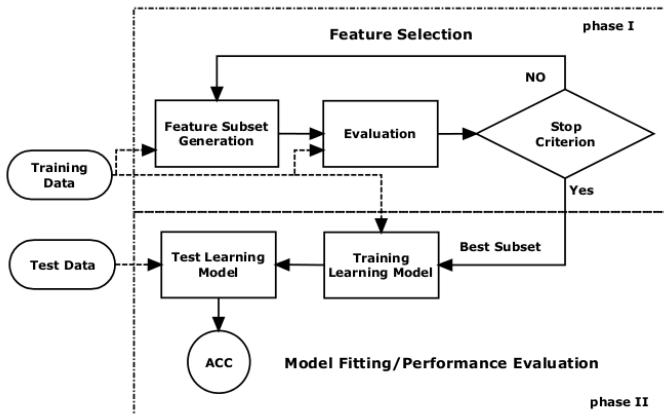
Feature redundancy

- Feature redundancy is usually presented in terms of feature correlation.
- Perfectly correlated features are truly redundant in the sense that no additional information is gained by adding them.
- Redundancy may exist between two uncorrelated features.
- Two highly correlated features may improve the accuracy \Rightarrow correlation cannot be adequately to feature redundancy.

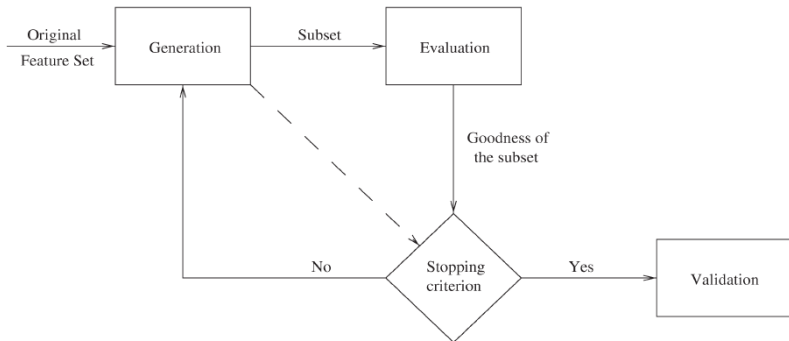
- Selecting the most relevant variables is usually suboptimal for building a predictor, particularly if the variables are redundant
- A subset of useful variables may exclude many redundant, but relevant, variables.
- Relevance does not imply optimality \equiv Let X_1, X_2, X_3 be binary features. Let the distribution of instances be uniform, and assume that the target concept is $J(X_1, X_2, X_3) = (X_1 \wedge X_2) \vee X_3$. In this case, all features are relevant but the optimal subset of features is $\{X_3\}$.



Feature selection workflow



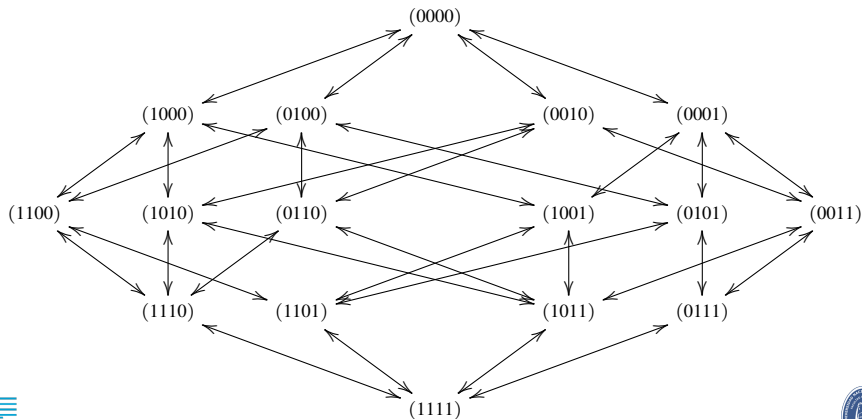
Steps in a typical feature selection method



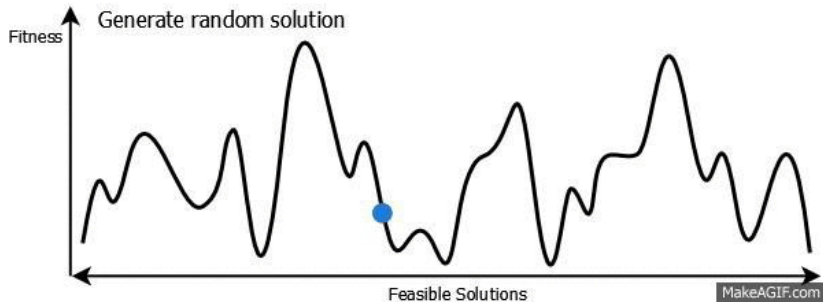
- Feature subset generation \equiv select subset candidate.
- Evaluation \equiv compute relevancy value of the subset.
- Stopping criterion \equiv determine whether subset is relevant.
- Validation \equiv verify subset validity.

Feature subset generation

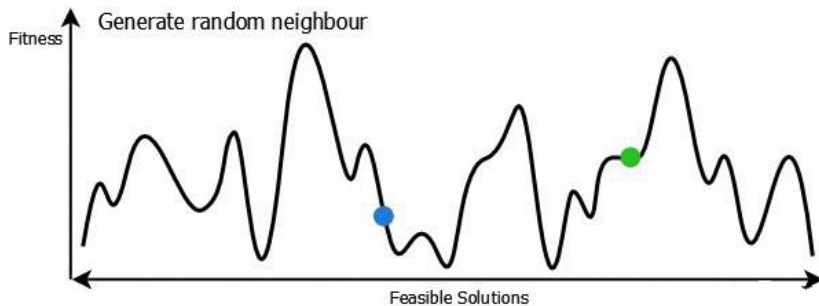
Search space



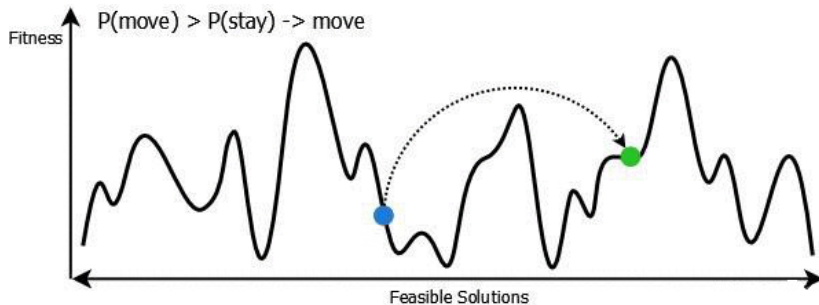
Search space



Search space



Search space



Feature subset generation

Approaches to examine the search space

- Complete \equiv it does a complete search for the optimal subset according to the evaluation function.
 - \Rightarrow Worst case: Exhaustive search ($\mathcal{O}(2^d)$).
 - \Rightarrow Optimality of the feature subset, according to the evaluation function, is guaranteed.
- Heuristic \equiv it generates the subsets under certain guidelines.
 - \Rightarrow Optimality is not guaranteed.
 - \Rightarrow Procedures very simple to implement and fast in producing results.
 - \Rightarrow Search space is usually quadratic ($\mathcal{O}(d^2)$).
 - Deterministic \equiv it generates the subsets in a predefined way.
 - Non deterministic \equiv it generates the subsets randomly.

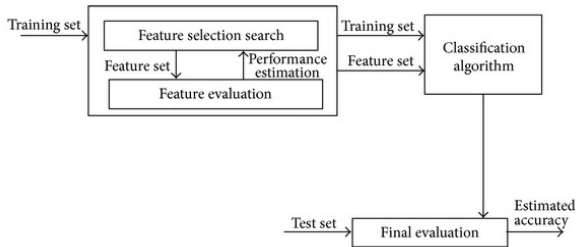
Evaluation

Determines the relevancy of the generated feature subset candidate towards the classification task.

Type of evaluation functions

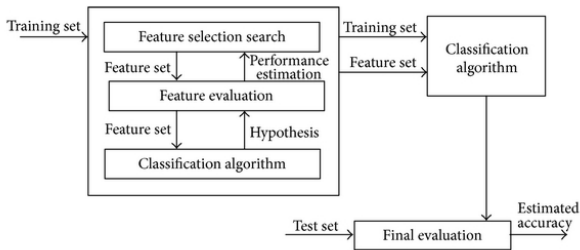
- Filter
 - 1 Distance (euclidean distance, Manhattan distance, etc.).
 - 2 Information (entropy, information gain, etc.)
 - 3 Dependency (correlation).
 - 4 Consistency (min-features bias).
- Wrapper (classifier).

Filter approach



- FS done as a preprocessing step.
- Subsets evaluated according to intrinsic properties of the data.
- Computationally fast \Rightarrow can be scaled to high-dimensional datasets.
- Drawback: Effect of FS on induction algorithm is not known.

Wrapper approach



- Wrappers use the learner as a black box to score the subsets of features according to their predictive power.
- The quality of feature subsets for classification is defined with respect to the induction algorithms.
- Drawback: Wrappers are slow.

Filters vs. wrappers

- Wrappers tend to have higher risk of overfitting than filters.
- Filter may lead to worse accuracy than wrappers.
- Filters are independent of the learner \Rightarrow FS done once for a given training dataset.

Categorization of feature selection methods

Search Evaluation	Complete		Heuristic	
			deterministic	non deterministic
Distance	<i>B&B</i> [31] <i>BFF</i> [42] <i>Seg84</i> [36] <i>EUBAFES</i> [34]	<i>Relief</i> [18] <i>ReliefF</i> [22]		
Information	<i>MDLM</i> [37]	<i>SFG</i> [4] <i>DT - CBL</i> [3] <i>DTM</i> [4] <i>KS96</i> [20] FCBF [43] <i>MIFS</i> [2] <i>CR</i> [40]	<i>PGVNS</i> [12]	
Dependence		<i>POE + ACC</i> [30] <i>PRESET</i> [28]		
Consistency	<i>FOCUS</i> [1] <i>Sch93</i> [35] <i>MIFES1</i> [32] <i>ABB</i> [24]	<i>SetCover</i> [6] <i>VCC</i> [41]	<i>LVF</i> [26] <i>SLV</i> [27] <i>QBB</i> [7]	
Error	<i>AMB&B</i> [11] <i>BS</i> [9] <i>LC</i> [14] <i>BC</i> [15] <i>PQSS</i> [9]	SFS [8] SBE [8] <i>SBE - SLASH</i> [5] <i>SFFS</i> [33] <i>BDS</i> [9] <i>RACE</i> [29] <i>RC</i> [10] <i>RACE</i> [29] <i>Oblivion</i> [23] <i>IS</i> [39] <i>RFE</i> [13]	<i>LVW</i> [25] <i>GA</i> [38] <i>SA</i> [9] <i>FSSEBNA</i> [16] SS [12]	

Sequential Forward Selection (SFS)

Main idea

- Deterministic heuristic search.
- Filter and wrapper approach.
- Complexity ($\mathcal{O}(d^2)$).
- Forward search.
- Starts with empty set.
- Each step, adds the best feature if its addition improves current solution.
- SFS performs best when the optimal subset is small.
- SFS is unable to remove features \Rightarrow the solution can get stuck in a local optimum.

Pseudocode

Procedure *Sequential Forward Search*

begin

1: $S \leftarrow \{\emptyset\}$

2: **repeat**

3: **foreach** $X_j \notin S$;

4: $J_j \leftarrow J(S \cup \{X_j\})$;

5: Let $j' \leftarrow \arg \max\{J_j\}$;

6: $S' \leftarrow S \cup \{X_{j'}\}$;

7: **if** $J(S') > J(S)$ **then**

8: $S \leftarrow S'$;

9: $J(S) \leftarrow J(S')$;

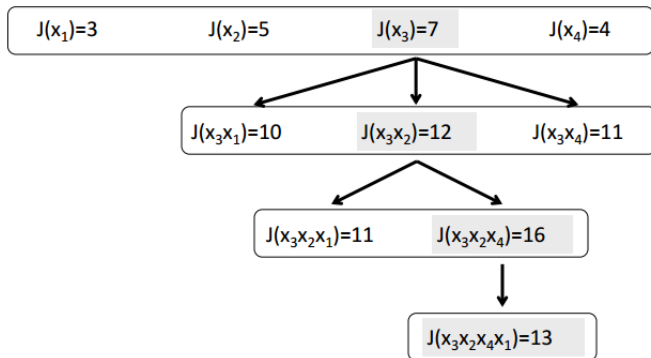
10: **until** $(J(S') \leq J(S) \parallel |S'| == d)$

end

Example of execution of SFS

Objective function

$$J(X) = -2x_1x_2 + 3x_1 + 5x_2 - 2x_1x_2x_3 + 7x_3 + 4x_4 - 2x_1x_2x_3x_4.$$



Sequential Backward Elimination (SBE)

Main idea

- Deterministic heuristic search.
- Filter and wrapper approach.
- Complexity ($\mathcal{O}(d^2)$).
- Starts with the full set set.
- Each step, removes the worst feature if its removal improves current solution.
- SBE performs best when the optimal subset is large.
- It is unable to reevaluate the usefulness of a feature after it has been discarded.

Pseudocode

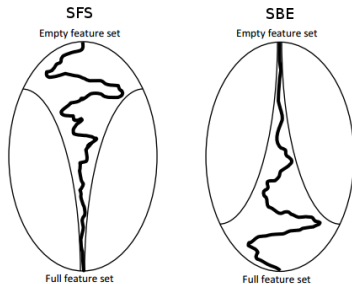
Procedure *Sequential Backward Elimination*

begin

```
1:  $S \leftarrow \{X_1, \dots, X_d\}$ 
2: repeat
3:   foreach  $X_j \in S$ ;
4:      $J_j \leftarrow J(S \setminus \{X_j\})$ ;
5:   Let  $j' \leftarrow \arg \max\{J_j\}$ ;
6:    $S' \leftarrow S \setminus \{X_{j'}\}$ ;
7:   if  $J(S') > J(S)$  then
8:      $S \leftarrow S'$ ;
9:      $J(S) \leftarrow J(S')$ ;
10: until  $(J(S') \leq J(S) \parallel |S'| == 1)$ 
```

end

SFS vs. SBE



- $|\mathcal{S}_{SFS}| \leq |\mathcal{S}_{SBE}|$.
- SFS may suffer of overfitting.
- $t_{SFS} \leq t_{SBE}$.
- SBE cannot be applied to medium high-dimensional data.
- Complexity ($\mathcal{O}(d^2)$) \Rightarrow not suitable for large high-dimensional data.

Fast Correlation Based Filter (FCBF)

Main idea

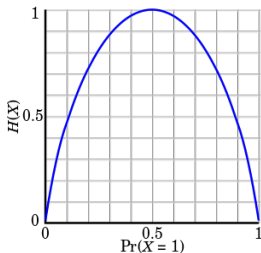
- Deterministic heuristic search.
- Filter approach \equiv information theory measures.
- Complexity:
 - best case: only one feature selected ($\mathcal{O}(d)$).
 - worst case: all features are selected ($\mathcal{O}(d^2)$).
- Two steps:
 - 1 Analysis of relevance.
 - 2 Analysis of redundance.

Definitions

Entropy

It measures the uncertainty about the value of a random variable X .

$$H(X) = - \sum_i P(x_i) \log_2(P(x_i)).$$



- Feature X with values $\{0, 1\}$.
- Entropy is 0 if there is no uncertainty.

Definitions

Conditional entropy

It measures the uncertainty about the value of X given the value of Y .

$$H(X|Y) = - \sum_j P(y_j) \sum_i P(x_i|y_j) \log_2(P(x_i|y_j)).$$

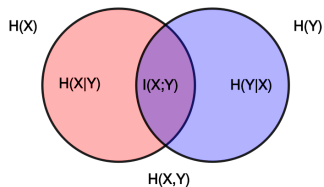
- $H(X|Y) = 0$ iff the value of X is completely determined by the value of Y .
- $H(X|Y) = H(X)$ iff X and Y are independent.

Definitions

Information Gain

It measures the reduction in uncertainty about the value of X given the value of Y

$$IG(X; Y) = H(X) - H(X|Y).$$

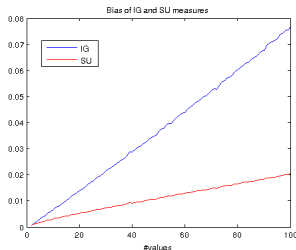


- $H(X) \equiv$ circle on the left (red and violet).
- $H(Y) \equiv$ circle on the right (blue and violet).
- $H(X, Y) \equiv$ area contained by both circles.
- $H(X|Y) \equiv$ red.
- $H(Y|X) \equiv$ blue.
- $I(X; Y) \equiv$ violet.

Definitions

Symmetrical Uncertainty (SU)

$$SU(X, Y) = 2 \left[\frac{IG(X; Y)}{H(X) + H(Y)} \right].$$



- 1000 examples generated randomly.
- 99 features with:
 - $X_1 = \{0, 1\}$,
 - $X_2 = \{0, 1, 2\}$,
 - \dots ,
 - $X_{99} = \{0, 1, \dots, 99\}$.
- Target class generated randomly \Rightarrow MI and SU values should be close to 0,



Definitions

Approximate Markov blanket (AMb)

Given two features X_i and X_j ($i \neq j$) so that $SU(X_j, \mathcal{Y}) \geq SU(X_i, \mathcal{Y})$, then X_j forms an approximate Markov blanket for X_i iff $SU(X_i, X_j) \geq SU(X_i, \mathcal{Y})$.

Predominant feature

Given a set of features S , a relevant feature is a predominant feature iff it does not have any AMb in S .

Analysis of relevance

- Relevance measure \equiv Symmetrical Uncertainty
 $SU(X_j, Y), j = 1, \dots, d.$
- Given δ , a feature X_j is irrelevant if $SU(X_j, Y) \leq \delta.$

Analysis of redundance

Markov blanket [21] \equiv Given a feature $X_j, M_j \subset X (X_j \notin M_j)$ is said to be a Markov blanket for X_j iff

$$P(X - M_j - \{X_j\}, \mathcal{Y} | X_j, M_j) = P(X - M_j - \{X_j\}, \mathcal{Y} | M_j).$$

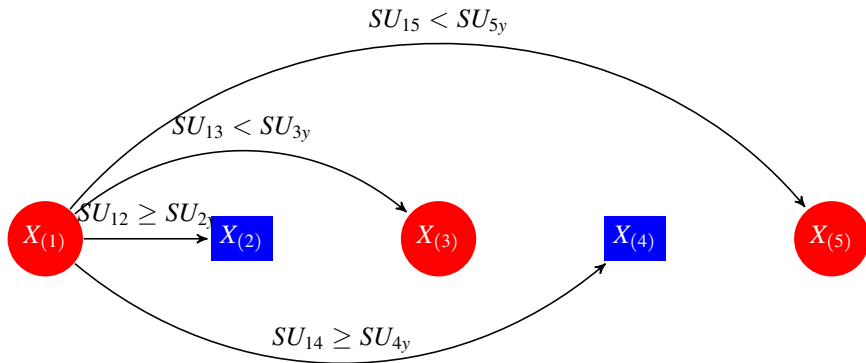
- M subsumes not only the information that X_i has about \mathcal{Y} but also about all of the other features.
- A feature $X_j \in S$ is redundant and, so, it can be removed from S if we find a Markov blanket M for X_j within $S.$

Analysis of redundance

- Features are ordered in descending order according to the SU values.
- First feature $X_{(1)}$ is a *predominant feature*.
- Second iteration \Rightarrow remove those features $X_{(j)}$ for which $X_{(1)}$ is an AMb.
- Second iteration \Rightarrow select next non-removed feature as *predominant feature* and remove those features for which, this feature forms an AMb.
- So on.

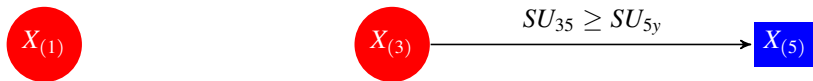
Example of analysis of redundance

Iteration 1



Example of analysis of redundance

Iteration 2



Pseudocode

Procedure *Fast Correlation Based Filter*

begin

```

1: for  $i = 1$  to  $d$  do
3:   calculate  $SU_{ic}$  for  $X_i$ ;
4:   if ( $SU_{ic} > \delta$ )
5:     append  $X_i$  to  $S'_{list}$ ;
6:   end;
7:   order  $S'_{list}$  in descending  $SU_{ic}$  value;
8:    $X_j = \text{getFirstElement}(S'_{list})$ ;
9:   do begin;
10:     $X_i = \text{getNextElement}(S'_{list}, X_j)$ ;
11:    if( $X_i \neq NULL$ )
12:      do begin;
13:        if( $SU_{ij} \geq SU_{ic}$ );
14:        remove  $X_i$  from  $S'_{list}$ ;
15:         $X_j = \text{getNextElement}(S'_{list}, X_j)$ ;
16:      end until ( $X_i == NULL$ );
17:       $X_j = \text{getNextElement}(S'_{list}, F_j)$ ;
18:    end until ( $X_j == NULL$ );
19:   $S_{best} = S'_{list}$ ;

```

end



Scatter Search (SS)

Main idea

- Heuristic and non deterministic method.
- Population based strategy.
- Evolution based on intensification and diversification strategies.

SS pseudocode

Procedure *Scatter search*

begin

1: *GeneratePopulation (InitPop);*

2: *GenerateReferenceSet (RefSet);*

3: **repeat**

4: **repeat**

5: *SelectSubset (Subset);*

6: *CombinationMethod (Subset, CurSol);*

7: *ImprovementMethod (CurSol, ImpSol);*

8: **until** (*StoppingCriterion₁*)

9: *UpdateReferenceSet (RefSet);*

10: **until** (*StoppingCriterion₂*)

end

Generate initial population

- Let L be an ordered subset with features of the subconjunto formado por los atributos con mayor poder predictivo (tal que $J(\{x_j\}) \geq J(\{x_{j+1}\})$).

```
1: Procedure Generate initial population
2: {
3:    $S \leftarrow \emptyset$ ;
4:   Order  $\{X_j\}, j = 1, \dots, d$  such that  $f(X_j) \geq f(X_{j+1})$ ;
5:    $L \leftarrow \{X_j\}, j = 1, \dots, k$  such that  $k \leq d$ ;
6:   repeat
7:     Select randomly  $X_{j^*} \in L$ ;
8:     if  $J(\{X_{j^*}\} \cup S) \geq J(S)$ 
9:        $S \leftarrow S \cup \{X_{j^*}\}$ 
10:     $L \leftarrow (L \setminus \{X_{j^*}\}) \cup \{X_j\}, X_j \notin L$ 
11:   until  $J(\{X_{j^*}\} \cup S) < J(S)$ 
12: }
```

Update the Reference Set

- $|RefSet| = |RefSet1| + |RefSet2|$.
- $RefSet1 \equiv$ quality and $RefSet2 \equiv$ diversity.

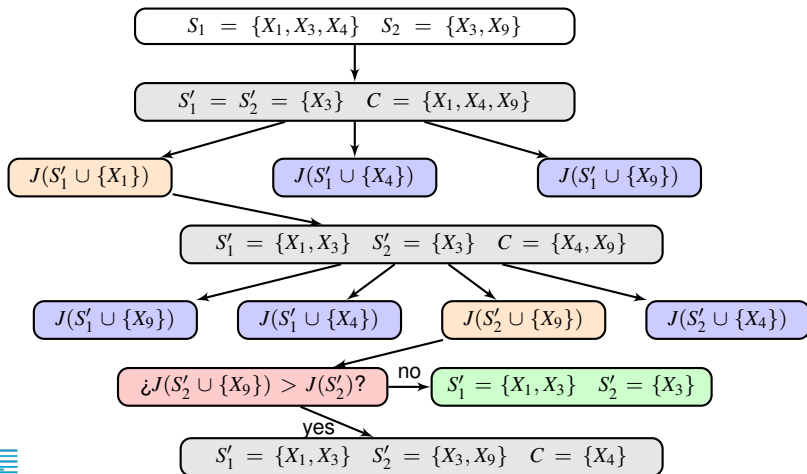
```
1: Procedure Update the Reference Set
2: {
3:    $RefSet \leftarrow \emptyset$ 
4:    $RefSet \leftarrow$  best solutions from Pop.
5:   Let  $C = \cup_{X_j \in RefSet} X_j$ 
6:   repeat  $\forall S \notin RefSet$ 
7:     Calculate  $Div(S, C) = |(S \cup C) \setminus (S \cap C)|$ 
8:     Let  $S^* \leftarrow \arg \max Div(S, C) : S \notin RefSet$ .
9:      $RefSet \leftarrow RefSet \cup S^*$ .
10:     $|RefSet| \leftarrow |RefSet| + 1$ .
11:   until  $|RefSet| = |RefSet1| + |RefSet2|$ 
12: }
```

Combination method

- $S_i \equiv$ Solution i .
- $S'_i \equiv$ new solution generated i .
- $(S_1, S_2) \rightarrow (S'_1, S'_2)$.

```
1: Procedure Greedy Combination
2: {
3:  $S'_1 = S'_2 \leftarrow S_1 \cap S_2, C = (S_1 \cup S_2) \setminus (S_1 \cap S_2)$ .
4:  $S'_1 \leftarrow S'_1 \cup \{X_{j^*}\} : X_{j^*} = \max_j \{J(S'_1 \cup \{X_j\})\}$ .
5: repeat
6:    $j_k^* : J(S'_k \cup \{X_{j_k^*}\}) = \max_j \{J(S'_k \cup \{X_j\})\}, k = 1, 2;$ 
7:   Let  $j^{**} = \max_k \{J(S'_k \cup \{X_{j_k^*}\})\}$ 
8:   if  $J(S'_k \cup \{X_{j^{**}}\}) > J(S'_k)$ 
9:      $S'_k \leftarrow S'_k \cup \{X_{j^{**}}\}$ 
10:     $C \leftarrow C \setminus \{X_{j^{**}}\}$ 
11: until  $J(S'_k \cup \{X_{j^{**}}\}) \leq J(S'_k), k = 1, 2$ 
12: }
```

Example of the combination method



Improvement method

- Let $CA = \{X_j : X_j \notin S\}$, ordered according to the evaluation method ($J(\{x_j\}) \geq J(\{x_{j+1}\})$).

```
1: Procedure Improvement method
2: {
3:    $j \leftarrow 0$ 
4:   repeat
5:      $j \leftarrow j + 1$ ;
6:     if  $J(S \cup \{X_j\}) \geq J(S)$ 
7:        $S \leftarrow S \cup \{X_j\}$ 
8:   until  $j \leftarrow |CA|$ 
9: }
```

More topics related to feature selection

- Stability of the FS strategies.
- FS applied to regression and clustering.
- Causal Feature Selection.

- [1] H. Almuallim and T. G. Dietterich.
Learning with many irrelevant features.
In A. P. . T. M. Press, editor, *Proceedings of the Ninth National Conference on Artificial Intelligence*, pages 547–552, 1992.
- [2] R. Battiti.
Using mutual information for selecting features in supervised neural net learning.
IEEE Transactions on Neural Networks, 5(4):537–550, 1994.
- [3] C. Cardie and N. Howe.
Empirical methods in information extraction.
In D. Fischer, editor, *Proceedings of the Fourteenth International Conference on Machine Learning*, pages 65–79. Morgan Kaufmann, 1997.

- [4] C. Cardie93.
Using decision trees to improve case-based learning.
In Proceedings of the Tenth International Conference on Machine Learning, pages 25–32, 1993.
- [5] R. Caruana and D. Freitag.
Greedy attribute selection.
In Proceedings of International Conference on Machine Learning, pages 28–36. AAAI Press / The MIT Press, 1994.
- [6] M. Dash.
Feature selection via set cover.
In Proceedings of IEEE Knowledge and Data Engineering Exchange Workshop, pages 165–171, 1997.
- [7] M. Dash and H. Liu.
Hybrid search of feature subsets.
In Pacific Rim International Conference on Artificial Intelligence, pages 238–249, 1998.

- [8] P. A. Devijver and J. Kittler.
Pattern Recognition: A Statistical Approach.
Prentice Hall International, 1982.
- [9] J. Doak.
An evaluation of feature selection methods and their application
to computer security.
Technical report, University of California, Department of
Computer Science, 1992.
- [10] P. Domingos.
Contextsensitive feature selection for lazy learners.
Artificial Intelligence Review, 11:227–253, 1997.
- [11] I. Foroutan and J. Sklansky.
Feature selection for automatic classification of non-gaussian
data.
IEEE Transactions on Systems, Man and Cybernetics,
17(2):187–198, 1987.

- [12] F. C. García-López, M. García-Torres, B. Melián-Batista, J. A. M. Pérez, and J. M. Moreno-Vega.
Solving the feature selection problem by a parallel scatter search.
European Journal of Operations Research, 169(2):477–489, 2006.
- [13] I. Guyon and A. Elisseeff.
An introduction to variable and feature selection.
Journal of Machine Learning Research, 3:1157–1182, 2003.
- [14] M. Ichino and J. Sklansky.
Feature selection for linear classifier.
In Proceedings of the Seventh International Conference on Pattern Recognition, pages 124–127. Morgan Kaufmann, 1984.

- [15] M. Ichino and J. Sklansky.
Optimum feature selection by zero-one programming.
IEEE Transactions on Systems, Man and Cybernetics,
14(5):737–746, 1984.
- [16] I. Inza, P. Larrañaga, R. Etxeberria, and B. Sierra.
Feature subset selection by bayesian networks based
optimization.
Artificial Intelligence, 123(1-2):157–184, 2000.
- [17] J. M. Keynes.
A treatise on probability.
Mcmillan and Co., 1921.
- [18] K. Kira and L. A. Rendell.
The feature selection problem.
In *In Proceedings of the Tenth National Conference on Artificial
Intelligence*, pages 129–134. Menlo Park: AAAI Press / The MIT
Press, 1992.

- [19] R. Kohavi and G. H. John.
Wrappers for feature subset selection.
Artificial Intelligence, 97(1-2):273–324, 1997.
- [20] D. Koller and M. Sahami.
Toward optimal feature selection.
In *Machine Learning: Proceedings of the Thirteenth International Conference*, pages 284–292, 1996.
- [21] D. Koller and M. Sahami.
Toward optimal feature selection.
In *Proceedings of the Thirteenth International Conference on Machine Learning*, pages 284–292, 1996.
- [22] I. Kononenko.
Estimating attributes: Analysis and extension of relief.
In *In Proceedings of the European Conference on Machine Learning*, pages 171–182. Springer-Verlag, 1994.

- [23] P. Langley and S. Sage.
Oblivious decision trees and abstract cases.
In Working Notes of the AAAI-94 Workshop on Case-Based Reasoning. AAAI Press, 1994.
- [24] H. Liu, H. Motoda, and M. Dash.
A monotonic measure for optimal feature selection.
In C. Nedellec and C. Rouveirol, editors, Machine Learning: ECML-98, pages 101–106. Springer-Verlag, 1998.
- [25] H. Liu and R. Setiono.
Feature selection and classification - a probabilistic wrapper approach.
In Proceedings of the Ninth International Conference on Industrial and Engineering Applications of AI and ES, pages 419–424, 1996.

- [26] H. Liu and R. Setiono.
A probabilistic approach to feature selection - a filter solution.
In *Proceedings of International on Machine Learning*, pages 319–327. Morgan Kaufmann, 1996.
- [27] H. Liu and R. Setiono.
Scalable feature selection for large sized databases.
In *Proceedings of the Fourth World Congress on Expert Systems*, volume 1. Morgan Kaufmann, 1998.
- [28] M. Modrzejewski.
Feature selection using rough sets theory.
In *Proceedings of the European Conference on Machine Learning*, pages 213–226, 1993.
- [29] A. W. Moore and M. S. Lee.
Efficient algorithms for minimizing cross validation error.
In *International Conference on Machine Learning*, pages 190–198, 1994.

- [30] A. N. Mucciardi and E. E. Gose.
A comparison of seven techniques for choosing subsets of pattern recognition.
IEEE Transactions on Computers, 20:1023–1031, 1971.
- [31] P. M. Narendra and K. Fukunaga.
A branch and bound algorithm for feature subset selection.
IEEE Trans. on Computer, 26(9):917–922, 1977.
- [32] A. L. Oliveira and A. S. Vincentelli.
Constructive induction using a non-greedy strategy for feature selection.
In Proceedings of the Ninth International Conference on Machine Learning, pages 355–360. Morgan Kaufmann, 1992.
- [33] P. Pudil, J. Novovicova, and J. Kittler.
Floating search methods in feature selection.
Pattern Recognition Letters, 15:1119–1125, 1994.

- [34] M. Scherf and W. Brauer.
Improving rbf networks by the feature selection approach
eubafes.
*In Proceedings of the 7th International Conference on Artificial
Neural Networks*, pages 391–396, 1997.
- [35] J. C. Schlimmer.
Efficiently inducing determinations: A complete and systematic
search algorithm that uses optimal pruning.
*In In Proceedings og the Tenth Inetrnational Conference on
Machine Learning*, pages 284–290. ICML93, 1993.
- [36] J. Segen.
Feature selection and constructive inference.
*In In Proceedings of the Seventh International Conference on
Pattern Recognition*, pages 1344–1346, 1984.

- [37] J. Sheinvald, B. Dom, and W. Niblack.
A modelling approach to feature selection.
In In Proceedings of the Tenth International Conference on Pattern Recognition, pages 535–539, 1990.
- [38] W. Siedlecki and J. Sklansky.
A note on genetic algorithms for large-scale feature selection.
Pattern Recognition Letters, 10:335–347, 1989.
- [39] H. Vafaie and I. F. Imam.
Feature selection methods: Genetic algorithms vs. greedy-like search.
In Proceedings of the International Conference on Fuzzy and Intelligent Control Systems, 1994.

- [40] H. Wang, D. Bell, and F. Murtagh.
Axiomatic approach to feature subset selection based on relevance.
IEEE Transactions on Pattern Analysis and machine Intelligence, 21(3):271–277, 1999.
- [41] K. Wang and S. Sundaresh.
Selecting features by vertical compactness of data.
In D. Heckerman, H. Mannila, and D. Pregibon, editors,
Proceedings of the Third International Conference on Knowledge Discovery and Data Mining, pages 275–278. AAAI Press, 1997.
- [42] L. Xu, P. Yan, and T. Chang.
Best first strategy for feature selection.
In *In Proceedings of the Ninth International Conference on Pattern Recognition*, pages 706–708, 1988.

- [43] L. Yu and H. Liu.
Efficient feature selection via analysis of relevance and
redundancy.
J. Mach. Learn. Res., 5:1205–1224, 2004.