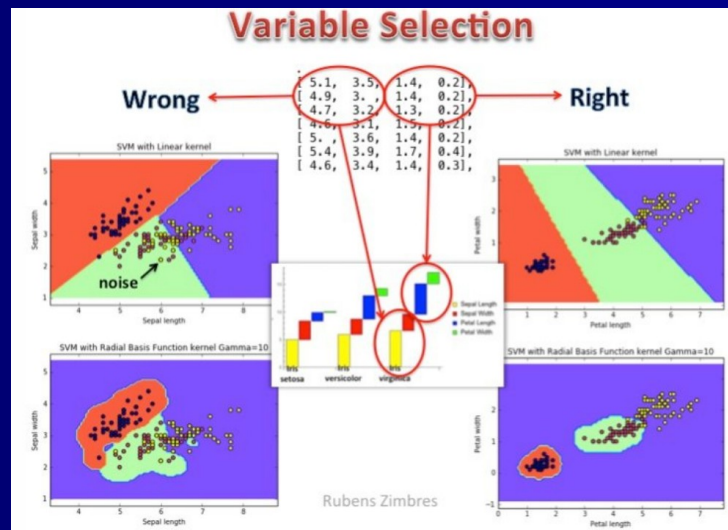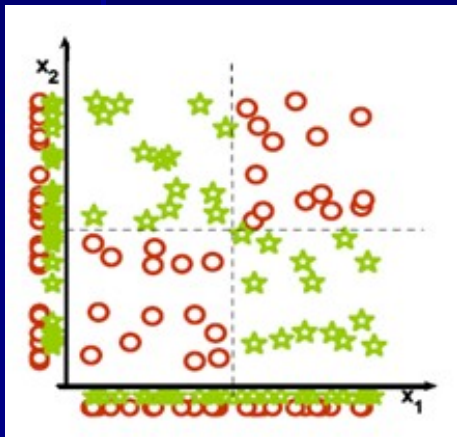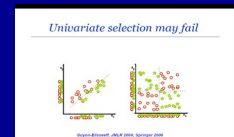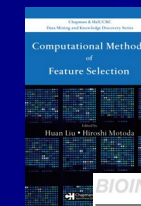# DIMENSIONALITY REDUCTION BY FEATURE SELECTION

**Iñaki Inza**
**Intelligent Systems Group, www.sc.ehu.es/isg**
**Computer Science Faculty**
**University of the Basque Country, Donostia - San Sebastian**

# OUTLINE

- The context: FS in supervised classification

- FS versus feature extraction

- Types of techniques

- Final remarks and ideas

- References and software

# FEATURE SELECTION (FS)
# FOR SUPERVISED CLASSIFICATION

- Fix the learning scenario: supervised classification



- Reduce the number of original features → $(X_1, X_2, ..., X_d)$
- Irrelevancy?
- Redundancy?

# FEATURE SELECTION (FS) FOR SUPERVISED CLASSIFICATION

- Fix the learning scenario: supervised classification



- Reduce the number of original features → $(X_1, X_2, …, X_d)$

- Improve accuracy
- Reduce costs
- Computational cost

# FEATURE EXTRACTION

- Feature selection ≠≠≠≠ Feature construction-extraction

- Feature extraction → PCA, SVD, PLS...

- Mathematical properties
- Intuition - interpretation

```
eigenvalue      proportion      cumulative
 2.91082          0.7277          0.7277          -0.581petallength-0.566petalwidth-0.522sepallength+0.263sepalwidth
 0.92122          0.23031         0.95801          0.926sepalwidth+0.372sepallength+0.065petalwidth+0.021petallength
```

```
Ranked attributes:
 0.2723  1 -0.581petallength-0.566petalwidth-0.522sepallength+0.263sepalwidth
 0.042   2 0.926sepalwidth+0.372sepallength+0.065petalwidth+0.021petallength
```

# FEATURE EXTRACTION
## PRINCIPAL COMPONENT ANALYSIS

# A MATURE FIELD
# APPLICATION DOMAINS

## Feature selection methods for text classification: a systematic literature review

Julliano Trindade Pintas[1] · Leandro A. F. Fernandes[1] · Ana Cristina Bicharra Garcia[2]

## An Extensive Empirical Study of Feature Selection Metrics for Text Classification

George Forman                          GFORMAN@HPL.HP.COM
Hewlett-Packard Labs
Palo Alto, CA, USA 94304

## Analysis of Local Appearance-based Face Recognition: Effects of Feature Selection and Feature Normalization

Hazım Kemal Ekenel, Rainer Stiefelhagen
Interactive Systems Labs, Computer Science Department, Universität Karlsruhe (TH)
Am Fasanengarten 5, 76131, Karlsruhe, Germany
{ekenel,stiefel}@ira.uka.de

Available online at www.sciencedirect.com

SCIENCE DIRECT®

Pattern Recognition Letters

Pattern Recognition Letters 25 (2004) 1377–1388

www.elsevier.com/locate/patrec

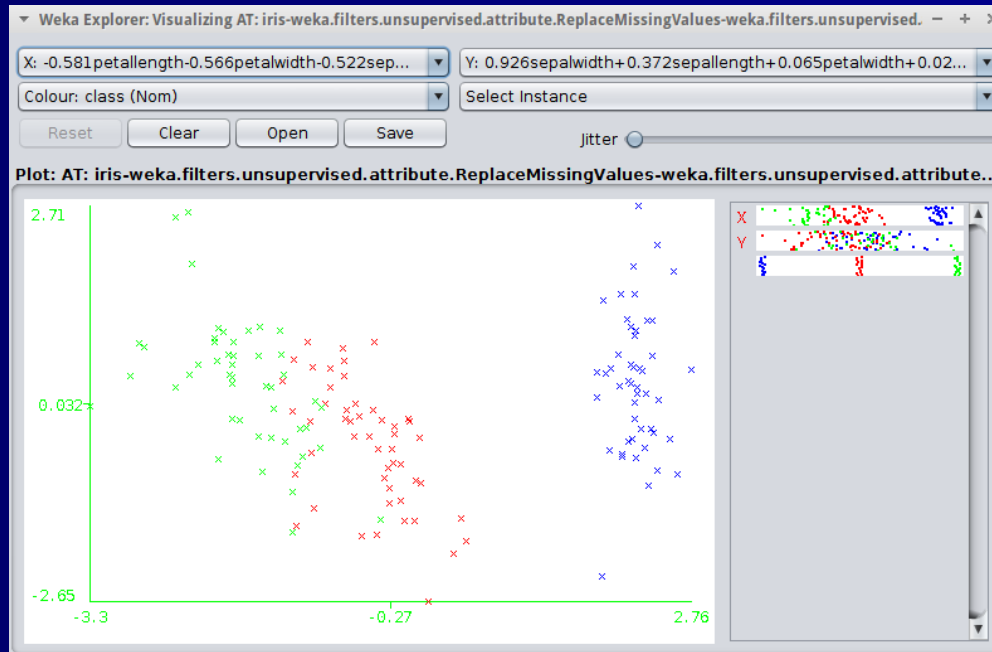## Feature selection in the independent component subspace for face recognition

H.K. Ekenel *, B. Sankur
Department of Electrical and Electronic Engineering, Bogazici University, Bebek 34342, Istanbul, Turkey

## MINIMUM REDUNDANCY FEATURE SELECTION FROM MICROARRAY GENE EXPRESSION DATA

CHRIS DING† and HANCHUAN PENG‡
†Computational Research Division and ‡Life Sciences/Genomics Division
Lawrence Berkeley National Laboratory, University of California
Berkeley, CA, 94720, USA

*BIOINFORMATICS*

REVIEW

*Gene expression*

## A review of feature selection techniques in bioinformatics

Yvan Saeys[1],*, Iñaki Inza[2] and Pedro Larrañaga[2]
[1]Department of Plant Systems Biology, VIB, B-9052 Ghent, Belgium and Bioinformatics and Evolutionary Genomics group, Department of Molecular Genetics, Ghent University, B-9052 Ghent, Belgium and [2]Department of Computer Science and Artificial Intelligence, Computer Science Faculty, University of the Basque Country, Paseo Manuel de Lardizabal 1, 20018 Donostia - San Sebastián, Spain

# TYPES OF FS TECHNIQUES

- Filter → selection independent of classifier
- Wrapper → ad-hoc features selected for a classifier


- Univariate filter


- Multivariate filter
- Multivariate wrapper


- Univariate → score individual features
- Multivariate → score feature subsets

# UNIVARIATE FILTER

- ORDER - RANK features
- By their correlation with the class-label

- Choose top-k features → learn classifier

- Correlation measures
- Chi-square, G², t-test, Fisher Criterion Scoring, Entropy-based metrics (mutual info)…

- How to choose $k$?
- Correlations among selected features?
- Complementarity?

```
Ranked attributes:
0.39065    5 odor
0.25795    8 gill-size
0.23312   12 stalk-surface-above-ring
0.21818   20 spore-print-color
0.20716   19 ring-type
0.19644    4 bruises?
0.19433   13 stalk-surface-below-ring
0.15815    7 gill-spacing
0.1376     9 gill-color
0.13106   14 stalk-color-above-ring
0.12204   15 stalk-color-below-ring
0.12137   17 veil-color
0.10081   21 population
0.09141   18 ring-number
0.08182    6 gill-attachment
0.06895   22 habitat
0.02952    1 cap-shape
0.02848   11 stalk-root
0.01815    2 cap-surface
0.01436    3 cap-color
0.00762   10 stalk-shape
0         16 veil-type
```

# UNIVARIATE FILTER CORRELATION by t-test



$$t_{test}(X_i, C) = \frac{|\mu_{X_i}^{c_1} - \mu_{X_i}^{c_2}|}{\sigma_{X_i}^{c_1} + \sigma_{X_i}^{c_2}}$$

# UNIVARIATE FILTER
# CORRELATION by MUTUAL-INFORMATION

- Based on the Entropy concept of a random variable

$$H(C) = -\sum_{c_i} p(c_i) \log_2 p(c_i)$$



- For example: information gain - mutual information, symmetrical uncertainty

$$MI(X_i, C) = IG(X_i, C) = H(C) - H(C|X_i)$$

$$SU(X_i, C) = \frac{2 \times MI(X_i, C)}{H(X_i) + H(C)}$$

# UNIVARIATE FILTER
## CORRELATION by MUTUAL-INFORMATION

# UNIVARIATE FILTER
## CORRELATION by CHI-SQUARE

The value of the test-statistic is

$$X^2 = \sum_{i=1}^{n} \frac{(O_i - E_i)^2}{E_i},$$

where

$X^2$ = Pearson's cumulative test statistic, which asymptotically approaches a $\chi^2$ distribution.
$O_i$ = an observed frequency;
$E_i$ = an expected (theoretical) frequency, asserted by the null hypothesis;
$n$ = the number of cells in the table.

|  | C=0 | C=1 |  |
|---|---|---|---|
| X=low | 16 | 2 | 18 |
| X=medium | 3 | 6 | 9 |
| X=high | 1 | 22 | 23 |
|  | 20 | 30 | 50 |

- If both variables were independent (predictor and class) → number of "expected" cases in each cell would be the product of the marginal, divided by the total number of cases:
  - e.g., "expected" number of cases in $cell_{11}$ is: (18x20)/50=7.2
  - while the "observed" number of cases in $cell_{11}$ is: 16
- Larger differences between expected and observed values in each cell → indicative of a "dependence-relationship" between predictor and class

# Univariate selection may fail



Guyon-Elisseeff, JMLR 2004; Springer 2006

# MULTIVARIATE FEATURE SELECTION
# <u>NUMBER</u> OF FEATURE SUBSETS



- Can the goodness of all feature subsets be evaluated? Hardly...
- $2^{10}=1024$; $2^{20}\approx1,048,576$; $2^{50}=1,125,899,906,842,624$



Kohavi-John, 1997

N features, $2^N$ possible feature subsets!

# HOW TO <u>EVALUATE</u> THE "GOODNESS" OF A FEATURE SUBSET?



Kohavi-John, 1997

N features, $2^N$ possible feature subsets!

# SUBSET EVALUATION (i) MULTIVARIATE FILTER SELECTION "CORRELATED FEATURE SELECTION (CFS)"

## Correlation-based Feature Selection for Discrete and Numeric Class Machine Learning

**Mark A. Hall**
MHALL@CS.WAIKATO.AC.NZ
Department of Computer Science, University of Waikato, Hamilton, New Zealand

$$Merit_s = \frac{k\overline{r_{cf}}}{\sqrt{k + k(k-1)\overline{r_{ff}}}} \quad (1)$$

where $Merit_S$ is the heuristic "merit" of a feature subset $S$ containing $k$ features, $\overline{r_{cf}}$ the average feature-class correlation, and $\overline{r_{ff}}$ the average feature-feature intercorrelation. Equation 1 is, in fact, Pearson's correlation, where all variables have been standardized. The numerator can be thought of as giving an indication of how predictive a group of features are; the denominator of how much redundancy there is among them. The heuristic handles irrelevant features as they will be poor predictors of the class. Redundant attributes are discriminated against as they will be

# SUBSET EVALUATION (i)
## MULTIVARIATE FILTER SELECTION
## "CORRELATED FEATURE SELECTION (CFS)"

- The $Merit_S$ function is calculated for each found feature subset *S*, for example:

$$Merit_S (X_3, X_6, X_8)?$$

- <u>Relevance</u> concept, $r_{cf}$: to be augmented
- Correlation of each {predictor ~ class} → enhances the merit-metric, e.g., corr($X_3$,*Class*), corr($X_6$,*Class*), corr($X_8$,*Class*)
- Irrelevant features → hurt e.g. K-NN

- <u>Redundancy</u> concept, $r_{ff}$: to be diminished
- Correlation among pairs of predictors → reduces the merit-metric, e.g., corr($X_3$, $X_6$), corr($X_6$, $X_8$), corr($X_8$, $X_3$)
- Redundant features → hurt e.g. naïve Bayes

- Any type of correlation measure can be used (t-test, MI...)!

$$Merit_s = \frac{k\overline{r_{cf}}}{\sqrt{k + k(k-1)\overline{r_{ff}}}}$$

# SUBSET EVALUATION (i)
## MULTIVARIATE FILTER SELECTION
## "MAX-RELEVANCE + MIN-REDUNDANCY"

Feature Selection Based on Mutual Information:
Criteria of Max-Dependency, Max-Relevance,
and Min-Redundancy

Hanchuan Peng, *Member*, *IEEE*, Fuhui Long, and Chris Ding

$$\max D(S, c), \quad D = \frac{1}{|S|} \sum_{x_i \in S} I(x_i; c)$$

$$\min R(S), \quad R = \frac{1}{|S|^2} \sum_{x_i, x_j \in S} I(x_i, x_j)$$

$$\max \Phi(D, R), \Phi = D - R$$

# MULTIVARIATE FILTER SELECTION ...MORE TECHNIQUES...

| Criterion | Full name | Authors |
|---|---|---|
| MIM | Mutual Information Maximisation | Lewis (1992) |
| MIFS | Mutual Information Feature Selection | Battiti (1994) |
| KS | Koller-Sahami metric | Koller and Sahami (1996) |
| JMI | Joint Mutual Information | Yang and Moody (1999) |
| MIFS-U | MIFS-'Uniform' | Kwak and Choi (2002) |
| IF | Informative Fragments | Vidal-Naquet and Ullman (2003) |
| FCBF | Fast Correlation Based Filter | Yu and Liu (2004) |
| AMIFS | Adaptive MIFS | Tesmer and Estevez (2004) |
| CMIM | Conditional Mutual Info Maximisation | Fleuret (2004) |
| MRMR | Max-Relevance Min-Redundancy | Peng et al. (2005) |
| ICAP | Interaction Capping | Jakulin (2005) |
| CIFE | Conditional Infomax Feature Extraction | Lin and Tang (2006) |
| DISR | Double Input Symmetrical Relevance | Meyer and Bontempi (2006) |
| MINRED | Minimum Redundancy | Duch (2006) |
| IGFS | Interaction Gain Feature Selection | El Akadi et al. (2008) |
| SOA | Second Order Approximation | Guo and Nixon (2009) |
| CMIFS | Conditional MIFS | Cheng et al. (2011) |

## Conditional Likelihood Maximisation: A Unifying Framework for Information Theoretic Feature Selection

**Gavin Brown**        GAVIN.BROWN@CS.MANCHESTER.AC.UK
**Adam Pocock**        ADAM.POCOCK@CS.MANCHESTER.AC.UK
**Ming-Jie Zhao**      MING-JIE.ZHAO@CS.MANCHESTER.AC.UK
**Mikel Luján**        MIKEL.LUJAN@CS.MANCHESTER.AC.UK
*School of Computer Science*
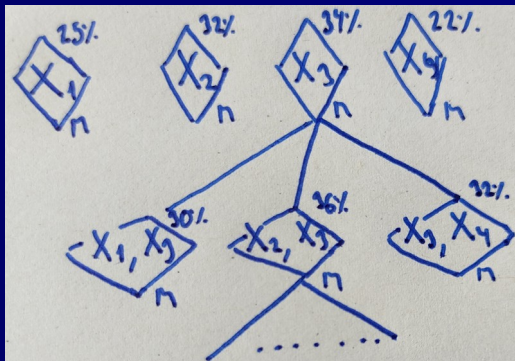*University of Manchester*
*Manchester M13 9PL, UK*

$$J_{mrmr}(X_k) = I(X_k;Y) - \frac{1}{|S|}\sum_{j \in S} I(X_k;X_j)$$

$$J'_{cmi}(X_k) = I(X_k;Y) - \sum_{j \in S} I(X_j;X_k) + \sum_{j \in S} I(X_j;X_k|Y)$$
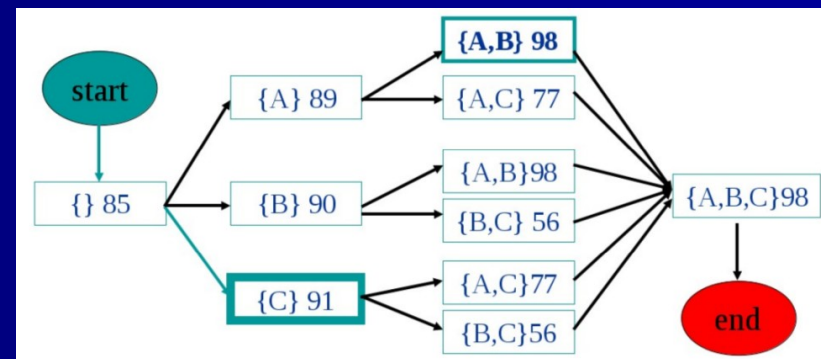
# SUBSET EVALUATION (ii)
# WRAPPER FEATURE SELECTION

- Fixed a <u>classification algorithm</u>

- Fixed a way to validate classifiers

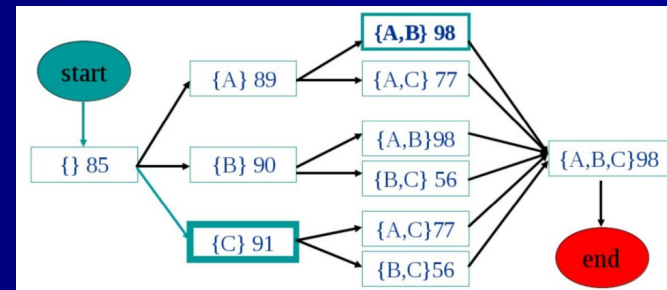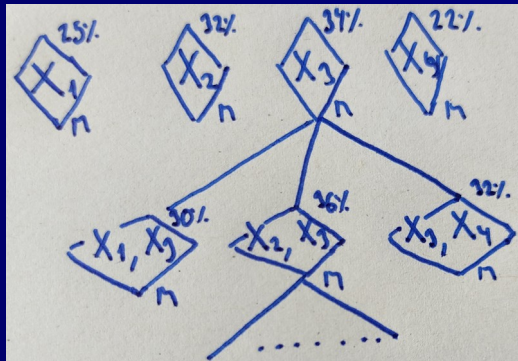- {Learn + Validate} classifier with the feature subset

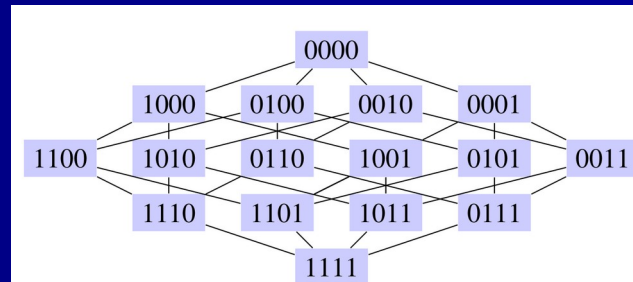| FEATURE SET | CLASSIFIER | PERFORMANCE |
|---|---|---|
| {A,B,C} | M | 98 % |
| {A,B} | M | 98 % |
| {A,C} | M | 77 % |
| {B,C} | M | 56 % |
| {A} | M | 89 % |
| {B} | M | 90 % |
| {C} | M | 91 % |
| {} | M | 85 % |

# SUBSET EVALUATION (ii) WRAPPER FEATURE SELECTION

- Fixed a classification algorithm
- Fixed a way to validate classifiers
- {Learn + Validate} classifier with the feature subset

- Be careful! CPU cost!
- Computational costs → naïve Bayes ‹‹‹ K-nearest neighbour ‹‹‹ neural networks…

- Accuracy estimation → k-fold cross-validation? → CPU cost!
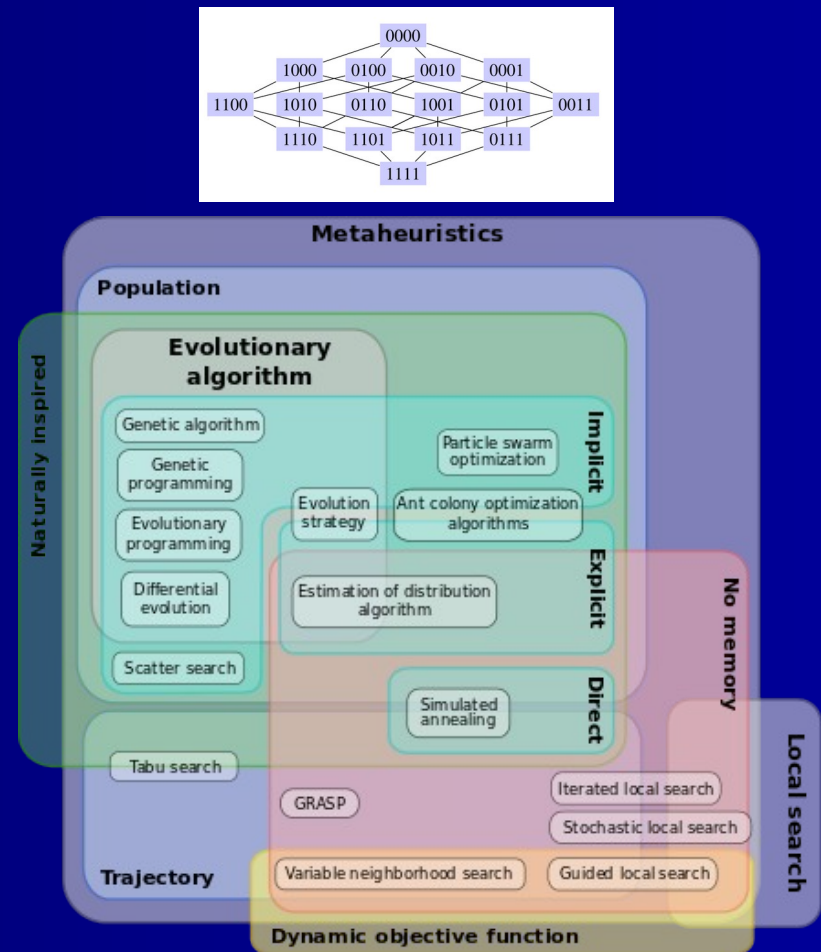
# MULTIVARIATE FEATURE SELECTION SEARCH PROBLEM

- Multivariate feature selection → search for the optimal feature subset

- How many different feature subsets? → $2^d$
- $2^{50}$=1,125,899,906,842,624 !!!!
- NP-hard problem

- Search heuristics → allowed
- Return → "suboptimal" solution

- Guarantee "optimal" solution → exhaustive search
- Computationally unfeasiable!

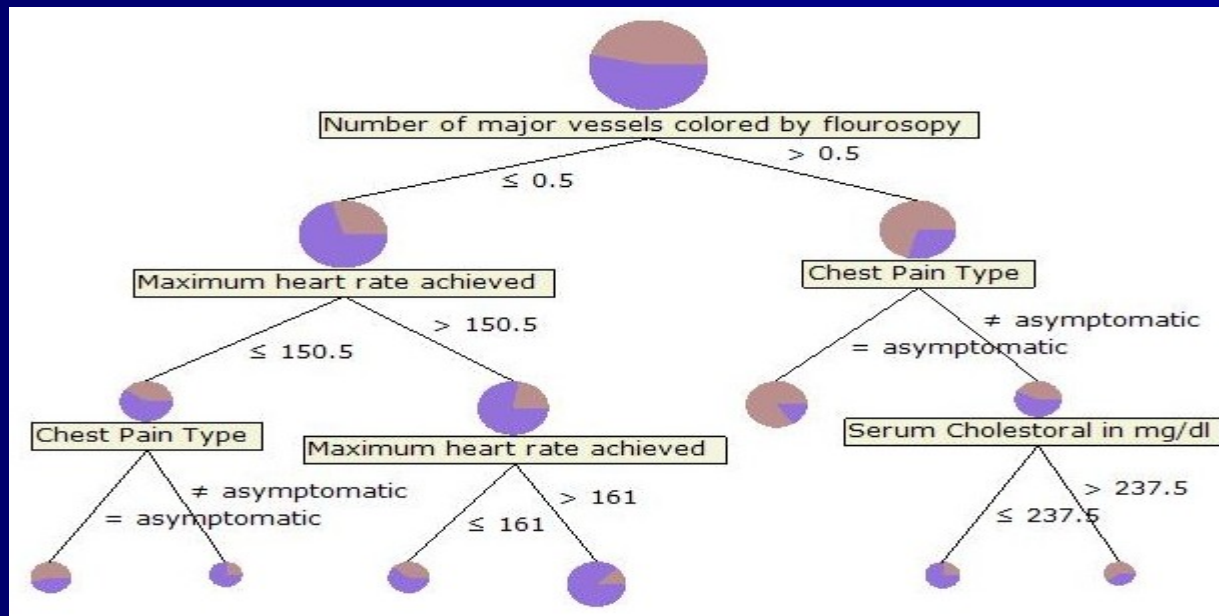- Common approach:
  - univariate filter + multivariate

# MULTIVARIATE FEATURE SELECTION SEARCH HEURISTICS

- Incremental – "local" heuristics
- Forward feature selection
- Backward feature elimination
- GRASP
- Simulated Annealing
- …

- "Population-based" - "Global"
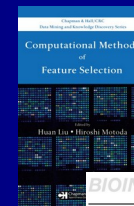- Genetic algorithms
- EDAs
- Ant-colony optimization
- …

# EMBEDDED FEATURE SELECTION

- Several classification algorithms
- "Embed" the capacity to discard initial features
- e.g. decision trees, random forest...

# REMARKS ON FS

- FS by Markov Blanket

- Honest model evaluation when applying FS

- Stability on FS

- MultiObjective FS

- FS in other learning scenarios?

- References and software

# MULTIVARIATE FILTER: MARKOV BLANKET

- Idea: selecting the features involved in the <u>Markov blanket of the class variable</u> in the learned <u>Bayesian network structure</u>

- Markov Blanket is a concept of Bayesian network theory
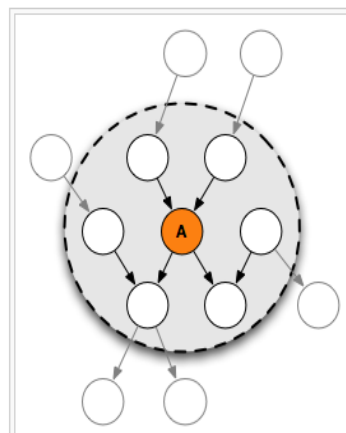- Constructed using an algorithm for Bayesian network learning

# HONEST ACCURACY ESTIMATION

## Overfitting in Making Comparisons Between Variable Selection Methods

**Juha Reunanen**
*ABB, Web Imaging Systems*
*P.O. Box 94, 00381 Helsinki, Finland*

JUHA.REUNANEN@FI.ABB.COM

*Data and text mining*

### Pitfalls of supervised feature selection

Pawel Smialowski[1,2,*], Dmitrij Frishman[1,2] and Stefan Kramer[3]

[1]Department of Genome Oriented Bioinformatics, Technische Universität München Wissenschaftszentrum Weihenstephan, Am Forum 1, 85350 Freising, [2]Helmholtz Zentrum Munich, National Research Center for Environment and Health, Institute for Bioinformatics, Ingolstädter Landstraße 1, 85764 Neuherberg and [3]Institut für Informatik/I12, Technische Universität München, Boltzmannstr. 3, 85748 Garching b. München, Germany

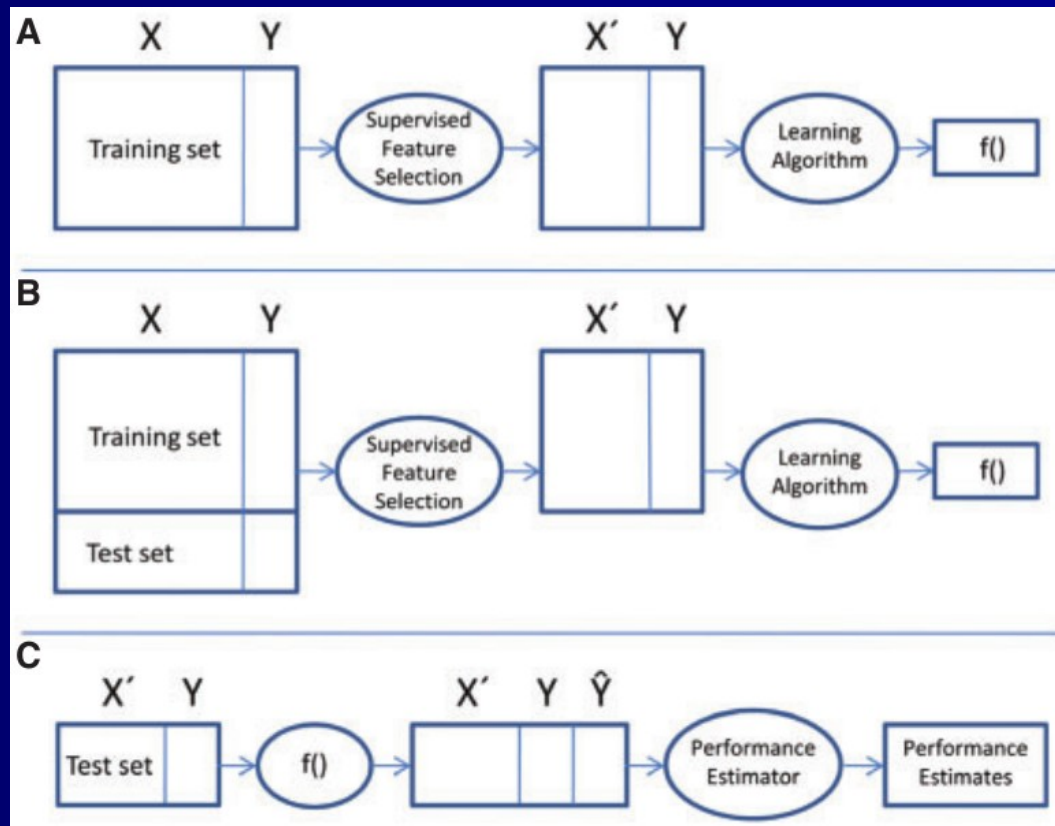### On feature selection protocols for very low-sample-size data

Ludmila I. Kuncheva[a], Juan J. Rodríguez[b,*]

[a] Bangor University, Dean Street, Bangor Gwynedd LL57 1UT, United Kingdom
[b] Universidad de Burgos, Escuela Politécnica Superior, Avda. de Cantabria s/n, Burgos 09006, Spain

# HONEST ACCURACY ESTIMATION



**Correct: A+C**          **Incorrect: B+C**

# STABILITY ON FEATURE SELECTION

- When applying different feature selection techniques
- Or using different seeds to randomly partition data

- Variablity on the subsets of selected features

- Proposing metrics to measure the *stability* on the subsets of selected features

- Learning a consensus subset?

## Measuring the Stability of Feature Selection

Authors                    Authors and affiliations

Sarah Nogueira ✉ , Gavin Brown

Proceedings of the Fifth IEEE International Conference on Data Mining (ICDM'05)

### Stability of Feature Selection Algorithms

Alexandros Kalousis, Julien Prados, Melanie Hilario
University of Geneva, Computer Science Department
Rue General Dufour 24, 1211 Geneva 4, Switzerland
{kalousis, prados, hilario}@cui.unige.ch

$$S_S(A, B) = 1 - \frac{|A| + |B| - 2|A \cap B|}{|A| + |B| - |A \cap B|} = \frac{|A \cap B|}{|A \cup B|}$$

$$I_C(A, B) = \frac{r - \frac{k^2}{n}}{k - \frac{k^2}{n}} = \frac{rn - k^2}{k(n - k)}$$

$$S_H(A, B) = 1 - \frac{|A \setminus B| + |B \setminus A|}{n}$$

# MULTIOBJECTIVE FEATURE SELECTION

- Optimization of 2 conflicting objectives:
  - accuracy
  - subset dimension

- Not a single solution
- Pareto Front: non-dominated solutions

- Why not other objective apart from accuracy?

# FS IN OTHER LEARNING SCENARIOS



Regression



Unsupervised learning



Semi-supervised learning



Positive-unlabeled learning

| $X_1$ | $X_2$ | ... | $X_n$ | $C_1$ | $C_2$ | ... | $C_m$ |
|-------|-------|-----|-------|-------|-------|-----|-------|
| $x_1^{(1)}$ | $x_2^{(1)}$ | ... | $x_n^{(1)}$ | $c_1^{(1)}$ | $c_2^{(1)}$ | ... | $c_m^{(1)}$ |
| $x_1^{(2)}$ | $x_2^{(2)}$ | ... | $x_n^{(2)}$ | $c_1^{(2)}$ | $c_2^{(2)}$ | ... | $c_m^{(2)}$ |
| ... | ... | ... | ... | ... | ... | ... | ... |
| $x_1^{(N)}$ | $x_2^{(N)}$ | ... | $x_n^{(N)}$ | $c_1^{(N)}$ | $c_2^{(N)}$ | ... | $c_m^{(N)}$ |

# FILTER vs. WRAPPER vs. EMBEDDED

| | Model search | | Advantages | Disadvantages | Examples |
|---|---|---|---|---|---|
| **Filter** | FS space → Classifier | Univariate | Fast<br>Scalable<br>Independent of the classifier | Ignores feature dependencies<br><br>Ignores interaction with the classifier | Chi-square<br>Euclidean distance<br>t-test<br>Information gain, Gain ratio [6] |
| | | Multivariate | Models feature dependencies<br>Independent of the classifier<br>Better computational complexity<br>than wrapper methods | Slower than univariate techniques<br>Less scalable than univariate<br>techniques<br>Ignores interaction with the classifier | Correlation based feature selection (CFS) [45]<br>Markov blanket filter (MBF) [62]<br>Fast correlation based<br>feature selection (FCBF) [136] |
| **Wrapper** | FS space / Hypothesis space / Classifier | Deterministic | Simple<br>Interacts with the classifier<br>Models feature dependencies<br>Less computationally intensive<br>than randomized methods | Risk of over fitting<br>More prone than randomized algorithms<br>to getting stuck in a local optimum<br>(greedy search)<br>Classifier dependent selection | Sequential forward selection (SFS) [60]<br>Sequential backward elimination (SBE) [60]<br>Plus $q$ take-away $r$ [33]<br>Beam search [106] |
| | | Randomized | Less prone to local optima<br>Interacts with the classifier<br>Models feature dependencies | Computationally intensive<br>Classifier dependent selection<br>Higher risk of overfitting<br>than deterministic algorithms | Simulated annealing<br>Randomized hill climbing [110]<br>Genetic algorithms [50]<br>Estimation of distribution algorithms [52] |
| **Embedded** | FS U Hypothesis space / Classifier | | Interacts with the classifier<br>Better computational complexity<br>than wrapper methods<br>Models feature dependencies | Classifier dependent selection | Decision trees<br>Weighted naive Bayes [28]<br>Feature selection using<br>the weight vector of SVM [44, 125] |

# REFERENCES

Chapman & Hall/CRC
Data Mining and Knowledge Discovery Series

**Computational Methods**
of
**Feature Selection**

Edited by
**Huan Liu • Hiroshi Motoda**

Chapman & Hall/CRC
Taylor & Francis Group

Journal of Machine Learning Research 3 (2003) 1157-1182          Submitted 11/02; Published 3/03

## An Introduction to Variable and Feature Selection

**Isabelle Guyon**                                              ISABELLE@CLOPINET.COM
*Clopinet*
*955 Creston Road*
*Berkeley, CA 94708-1501, USA*

**André Elisseeff**                                            ANDRE@TUEBINGEN.MPG.DE
*Empirical Inference for Machine Learning and Perception Department*
*Max Planck Institute for Biological Cybernetics*
*Spemannstrasse 38*
*72076 Tübingen, Germany*

Journal of Machine Learning Research 13 (2012) 27-66          Submitted 12/10; Revised 6/11; Published 1/12

## Conditional Likelihood Maximisation: A Unifying Framework for Information Theoretic Feature Selection

**Gavin Brown**                                    GAVIN.BROWN@CS.MANCHESTER.AC.UK
**Adam Pocock**                                    ADAM.POCOCK@CS.MANCHESTER.AC.UK
**Ming-Jie Zhao**                                  MING-JIE.ZHAO@CS.MANCHESTER.AC.UK
**Mikel Luján**                                    MIKEL.LUJAN@CS.MANCHESTER.AC.UK
*School of Computer Science*
*University of Manchester*
*Manchester M13 9PL, UK*

# SOFTWARE FOR FS