

RELEVANCE AND MUTUAL INFORMATION BASED FEATURE DISCRETIZATION

Artur J. Ferreira^{1,3}

Mário A. T. Figueiredo^{2,3}

¹Instituto Superior de Engenharia de Lisboa

²Instituto Superior Técnico

³Instituto de Telecomunicações

Lisboa, PORTUGAL

ICPRAM2013, Barcelona, Spain, 15-18 February, 2013

Outline

Introduction

Background

- Feature Discretization (FD)

Proposed Methods

- Relevance-Based Linde-Buzo-Gray (R-LBG)

- Mutual Information Discretization (MID)

Experimental Evaluation

- R-LBG and MID

- Unsupervised FD

- Supervised FD

Conclusions

Introduction and Motivation

Feature discretization (FD) in machine learning:

1. mandatory in some cases, optional in others
2. provides noise robustness, by ignoring (irrelevant) fluctuations on the data
3. yields compact data representations, reducing the storage requirements
4. may reduce training time and improve classification accuracy (can be seen as a form of regularization)
5. may be coupled with *feature selection* (FS), to further improve the the performance of some learning methods

Feature Discretization: taxonomy

FD techniques are usually categorized along five axes (Witten and Frank, 2005):

1. **unsupervised** vs **supervised**;
2. **static** (single pass assuming independent features) vs **dynamic** (taking dependencies into account);
3. **global** (discretizes the entire feature space) vs **local** (discretizes some features, as needed);
4. **top-down** (splitting) vs **bottom-up** (merging)
5. **direct** (sets a priori the number of bits per feature) vs **incremental**

Feature Discretization: quality indicators

The quality of discretization is usually assessed by two indicators:

- **generalization error**
- **complexity** (number of bits per feature instance)

Some facts from the FD literature:

- naturally, supervised FD may yield better classifiers (Dougherty et al., 1995; Witten and Frank, 2005)
- however, some unsupervised FD methods have been found to perform well on some types of data (Yang and Webb, 2001)
- no technique is uniformly better than all the others

Feature Discretization: unsupervised methods

Some commonly used **unsupervised** FD methods (Witten and Frank, 2005):

- *equal-interval binning* (EIB) - uniform quantization
- *equal-frequency binning* (EFB) - non-uniform quantization in which the number of occurrences in each interval is the same (Chiu et al., 1991)
- *proportional k-interval discretization* (PkID) - the number/size of intervals depend on the number of training instances (Yang and Webb, 2001)
- *unsupervised Linde-Buzo-Gray* (U-LBG) - discrete features with minimum *mean square error* (MSE) w.r.t. the original ones (Ferreira and Figueiredo, 2012)

U-LBG1/2 Algorithms (key ideas)

U-LBG1 and U-LBG2 are two versions of the U-LBG approach:

- rationale - low MSE between discrete and original features is adequate for learning
- the LBG algorithm is applied individually to each feature
- U-LBG1 uses a variable number of bits per feature, being stopped when:
 - the MSE distortion falls below some threshold Δ
 - or the maximum number of bits per feature q is reached
- U-LBG2 uses a fixed number of bits per feature, q

Feature Discretization: supervised methods

Some commonly used **supervised** FD methods:

- *information entropy minimization* (IEM) - uses a *minimum entropy criterion* in a top-down approach (Fayyad and Irani, 1993)
- *IEM variant* (IEMV) - uses MDL to control the number of different values intervals for each feature (Kononenko, 1995)
- *class-attribute interdependence maximization* (CAIM) - maximizes the class-attribute interdependence (Kurgan and Cios, 2004)
- *class-attribute contingency coefficient* (CACC) (Tsai et al., 2008)

Our Proposals for FD

In this paper, we propose two FD methods:

1. a *static, global, top-down, incremental*, relevance-based method for unsupervised or supervised learning
→ **Relevance-based LBG (R-LBG)**
2. a *static, global, top-down, incremental*, and *supervised* method based on the maximization of the *mutual information* (MI) between each feature and the class label
→ **Mutual Information Discretization (MID)**

Proposal 1: R-LBG Algorithm (key ideas)

The main characteristics of the R-LBG algorithm are as follows:

- applies the (unsupervised) LBG algorithm, with an incremental number of bits per feature
- it uses a (supervised or unsupervised) relevance function, $@rel$, and a (nonnegative) stopping factor ϵ
- $@rel$, producing non-negative values, is applied after each discretization
- for each feature, \tilde{X}_i , discretization is halted at b ($< q$) bits, whenever $@rel(\tilde{X}_i^{(b)}) - @rel(\tilde{X}_i^{(b-1)}) < \epsilon$
- setting $\epsilon = 0$, leads to the minimum number of bits that ensures maximum relevance
- it is a generalization of the U-LBG1 and U-LBG2 techniques

R-LBG Algorithm (relevance function)

Some choices for the unsupervised relevance function $@rel$:

- $@rel = MSE$ between original and discrete features, we have the unsupervised U-LBG1/2 approaches
- the quotient between the variance of the discrete feature and the number of discretization intervals

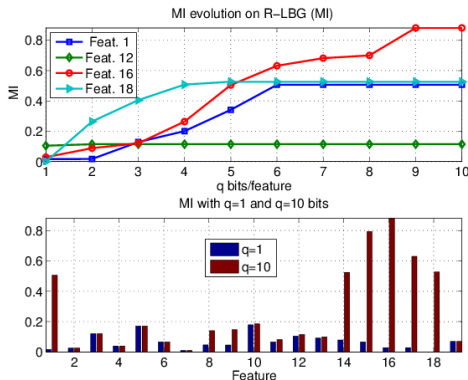
$$@rel(\tilde{X}_i^{(b)}) = NVAR(\tilde{X}_i) = \text{var}(\tilde{X}_i) / 2^b$$

For the supervised case, $@rel$ can be computed as:

- the MI between discretized features \tilde{X}_i and the class label \mathbf{y}
- the well-known Fisher's ratio (using the same operands)
- any ranking criterion used in feature selection

R-LBG: some insight on the relevance

R-LBG ($@rel = MI$) on the Hepatitis dataset ($d = 19$ features)



Top: MI as a function of the number of bits $q \in \{1, \dots, 10\}$, for features 1, 12 (which is categorical), 16, and 18.

Bottom: MI with $q = 1$ and $q = 10$ bits.

Proposal 2: Mutual Information Discretization

The key motivations for the MID algorithm are as follows:

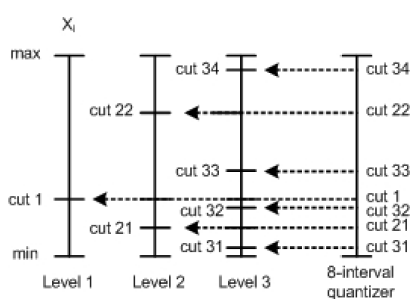
- good FS criteria should also be adequate for FD
- MI between features and class labels \mathbf{y} is adequate for FS
- the Hellman-Raviv (1970) and Santhi-Vardi (2006) bounds relate the Bayes error with the MI

$$err_{Bayes}(\tilde{X}_i) \leq \frac{1}{2} H(\mathbf{y}|\tilde{X}_i) \quad err_{Bayes}(\tilde{X}_i) \leq 1 - 2^{-H(\mathbf{y}|\tilde{X}_i)}.$$

- Recall that $MI(\tilde{X}_i; \mathbf{y}) = \underbrace{H(\mathbf{y})}_{\text{fixed}} - H(\mathbf{y}|\tilde{X}_i).$

MID: how does it work

- illustration for a given feature ($q = 3$ bit)
- it computes the cut-points that maximize the MI of the discrete feature with the class label



MID is an incremental and recursive partition algorithm

MID: fixed and variable versions

We propose two versions of MID:

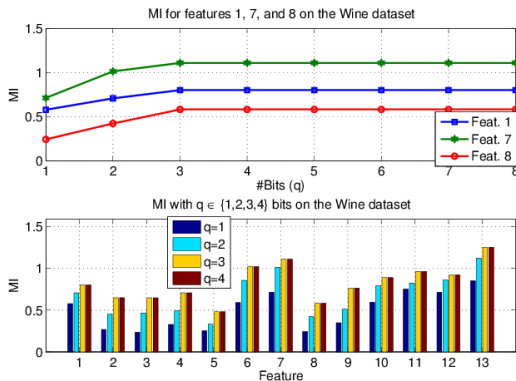
- MID-fixed, applies MID with q bits per feature
- MID-variable, allocates *up to* q bits per feature

MID-variable is controlled by a (nonnegative) stopping factor ϵ :

- $MI(\tilde{X}_i^{(b)}; \mathbf{y})$ is computed
- for each feature, discretization is halted at b bits, whenever $MI(\tilde{X}_i^{(b)}; \mathbf{y}) - MI(\tilde{X}_i^{(b-1)}; \mathbf{y}) < \epsilon$
- setting $\epsilon = 0$ leads to the minimum number of bits that maximizes the MI
- for a given q , MID-variable will produce fewer discretization intervals than MID-fixed

MID: evolution of MI

MI evolution as a function of the number of bits on the Wine dataset ($d = 13$ features)



Top: MI for features 1, 7, and 8, with $q \in \{1, \dots, 8\}$. Bottom: MI between discretized features and the class label, for $q \in \{1, 2, 3, 4\}$

Experimental Evaluation: Task and Datasets

- Supervised classification with linear SVM, naïve Bayes (NB), and k-nearest neighbors (KNN)
- 10-fold cross-validation
- UCI, microarray^{\$}, and face images[#] datasets with d features, c classes, and n patterns (in some cases, $d \gg n$)

Dataset	d	c	n	Dataset	d	c	n
Wine	13	3	178	Leukemia1 ^{\$}	5327	3	72
Hepatitis	19	2	155	TOX-171 ^{\$}	5748	4	171
Ionosphere	34	2	351	Brain-Tumor1 ^{\$}	5920	5	90
Colon ^{\$}	2000	2	62	ORL10P [#]	10304	10	100
SRBCT ^{\$}	2309	4	83	Prostate-Tumor ^{\$}	10509	2	102
AR10P [#]	2400	10	130	Leukemia2 ^{\$}	11225	3	72
PIE10P [#]	2420	10	210	GLI-85 ^{\$}	22283	2	85

Experimental Results: R-LBG and MID 1/3

R-LBG ($@rel = MI$) and *MID-variable* with $q = 4$ and linear SVM classifier

First row: total number of bits per instance

Second row: test error rate (%)

Dataset / No FD	R-LBG (MI)		MID variable	
	$\epsilon = 0$	$\epsilon = 0.1$	$\epsilon = 0$	$\epsilon = 0.1$
Wine	52.0	30.6	38.3	26.2
3.9	2.8	1.7	3.4	2.8
Hepatitis	46.5	68.8	28.5	65.6
21.3	15.5	21.9	18.7	18.1
Ionosphere	129.0	102.4	73.0	85.0
12.8	14.0	12.5	9.4	5.7
Colon	7954.6	7564.0	4682.0	6151.9
17.7	19.4	14.5	19.4	14.5
SRBCT	9222.5	8827.7	7144.2	7180.3
0.0	0.0	0.0	0.0	0.0
AR10P	9599.8	9583.2	8620.4	8640.4
0.8	0.8	0.8	0.0	0.0
PIE10P	9679.9	9662.5	8550.7	8543.4
0.0	0.0	0.0	0.0	0.0

Experimental Results: R-LBG and MID 2/3

R-LBG ($@rel = MI$) and *MID-variable* with $q = 4$ and linear SVM classifier

First row: total number of bits per instance

Second row: test error rate (%)

Dataset / No FD	R-LBG (MI)		MID variable	
	$\epsilon = 0$	$\epsilon = 0.1$	$\epsilon = 0$	$\epsilon = 0.1$
Leukemia1 8.3	21248.2	19818.7	14636.9	15555.9
	4.2	5.6	8.3	6.9
TOX-171 14.6	22988.5	21439.2	19012.4	20070.8
	2.3	2.9	4.1	4.1
Brain-Tumor1 11.1	23649.6	22174.4	17531.0	17436.5
	8.9	10.0	10.0	10.0
ORL10P 1.0	41215.6	41195.1	37410.3	37410.2
	1.0	1.0	2.0	2.0
Prostate-Tumor 10.8	41735.0	40431.3	25493.1	36598.8
	7.8	7.8	7.8	7.8
Leukemia2 5.6	44300.1	40072.9	31124.0	30255.4
	1.4	1.4	1.4	1.4
GLI-85 10.6	88561.7	84364.2	54906.9	72131.8
	8.2	8.2	8.2	8.2

Experimental Results: R-LBG and MID 3/3

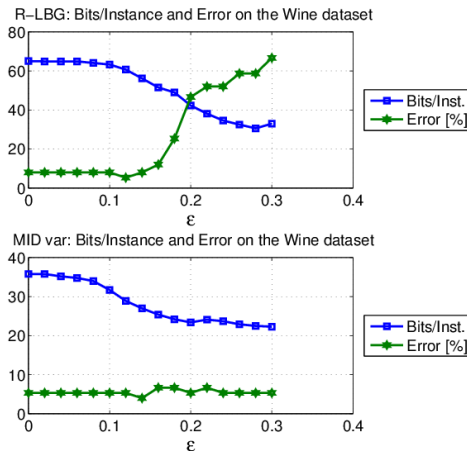
Some comments on these results:

- Obviously, R-LBG with $\epsilon = 0$ yields a larger number of bits per instance, as compared with $\epsilon = 0.1$
- R-LBG with $\epsilon = 0$ attains maximum relevance and better accuracy than with $\epsilon = 0.1$
- MID-variable with $\epsilon = 0$ yields the minimum number of bits that ensure the maximum MI
- MID-variable with $\epsilon = 0$ usually attains the better results with a few exceptions

The Friedman test reported a p-value of $0.04164 < 0.05$

R-LBG and MID: sensitivity to the ϵ parameter

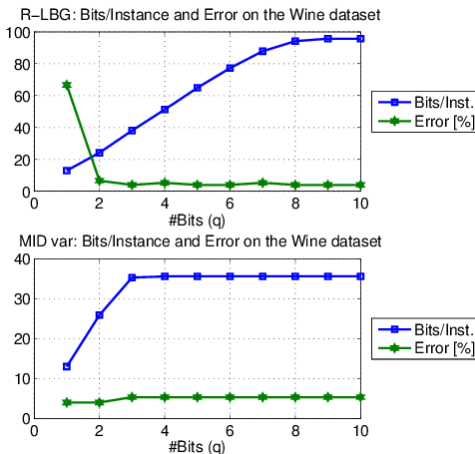
Average bits/instance and test set error rate (% , NB classifier) as function of the ϵ parameter



Wine dataset with $q = 5$ bits. Top: R-LBG (MI) Bottom: MID-variable

R-LBG and MID: sensitivity on the q parameter

Average bits/instance and test set error rate (% , NB classifier) as function of the q parameter



Wine dataset with $\epsilon = 0.05$. Top: R-LBG (MI) Bottom: MID-variable

Experimental Results: Unsupervised FD

Unsupervised FD with $q = 3$ bit/feature, $@rel = NVAR$, and $\epsilon = 0.25$

First row: total number of bits per instance

Second row: test error rate (%), linear SVM classifier

Dataset	No FD	Previous methods					Proposed
		EIB	EFB	PkID	U-LBG1	U-LBG2	R-LBG
AR10P		7200.0	7200.0	9267.8	7200.0	7200.0	6568.6
	0.8	0.0	0.8	0.8	0.8	0.8	1.5
PIE10P		7260.0	7260.0	9680.0	7260.0	7260.0	3774.4
	0.0	0.0	0.0	0.0	0.0	0.0	0.0
Leuk.1		15981.0	15981.0	15981.0	15981.0	15981.0	5733.0
	5.6	2.8	4.2	4.2	4.2	4.2	2.8
TOX-171		17244.0	17244.0	22992.0	17244.0	17244.0	5847.6
	9.9	1.2	1.8	1.2	1.8	1.8	8.2
B-Tumor1		17760.0	17760.0	23680.0	17760.0	17760.0	6085.4
	13.3	8.9	11.1	11.1	8.9	8.9	11.1
ORL10P		30912.0	30912.0	41216.0	30912.0	30912.0	19385.4
	1.0	1.0	1.0	1.0	1.0	1.0	1.0
P-Tumor		31527.0	31527.0	42035.4	31520.0	31527.0	11394.0
	10.8	8.8	8.8	8.8	8.8	8.8	8.8
Leuk.2		33675.0	33675.0	33675.0	33675.0	33675.0	12431.4
	4.2	2.8	2.8	2.8	2.8	2.8	4.2
GLI-85		66849.0	66849.0	66849.0	66849.0	66849.0	25118.7
	14.1	10.6	8.2	8.2	9.4	9.4	8.2

Experimental Results: Supervised FD

Supervised FD with $@rel = MI$ and $\epsilon = 0.1$

First row: total number of bits per instance

Second row: test error rate (%), linear SVM classifier

Dataset	Previous methods				Proposed methods		
	IEM	IEMV	CAIM	CACC	R-LBG	MIDf	MIDv
AR10P	12903.6	7138.4	7200.0	7200.0	7145.6	7200.0	7266.3
	2.3	20.0	0.8	0.0	0.0	0.0	0.0
PIE10P	9103.4	5264.0	7260.0	7260.0	7077.3	7260.0	7154.7
	0.0	1.9	0.0	0.0	0.0	0.0	0.0
Leuk.1	28435.3	26034.7	*	*	14656.1	15981.0	14278.0
	40.3	56.9	*	*	4.2	2.8	2.8
TOX-171	36134.8	28253.7	*	*	15725.6	17244.0	15824.2
	5.8	2.9	*	*	2.9	3.5	4.7
B.-Tumor1	32808.3	27133.5	*	*	15674.3	17760.0	16343.6
	20.0	35.6	*	*	11.1	10.0	8.9
ORL10P	26475.7	24176.8	*	*	30863.0	30912.0	30786.9
	9.0	1.0	*	*	2.0	2.0	2.0
P.-Tumor	54395.6	51964.7	*	*	30695.0	31527.0	28506.3
	12.7	11.8	*	*	5.9	6.9	7.8
Leuk.2	48380.1	40447.3	*	*	28857.3	33675.0	28670.8
	8.3	6.9	*	*	2.8	2.8	2.8
GLI-85	135866.9	130689.1	*	*	64065.0	66849.0	58633.4
	11.8	12.9	*	*	9.4	9.4	10.6

Conclusions

Our FD methods (Relevance-LBG and MID):

- for both unsupervised and supervised FD, attain equal or better results than state-of-the-art techniques
- scale well for high-dimensional data, contrary to other approaches
- have shown the adequacy of MI between features and class labels for discretization
- attain adequate discretizations with a variable number of bits per feature

In future work, we will explore the embedding of feature selection in the discretization process