



Data Mining & Knowledge Discovery

Class 6 – Feature Selection & Dimensionality Reduction 2025

Prof. Heitor Silvério Lopes Prof. Thiago H. Silva

Tasks x Methods in Data Mining

Tasks	Methods
Classification	Decision trees (C4.5), Classification rules, k-nearest-neighboors, Random forest, Support vector machine, Bayesian classifier, Neural network, Adaboost
Association Rules	Apriori, FP-growth, Eclat, Zigzag
Regression	Linear Regression, Polynomial regression, Logistic regression
Feature Selection & Dimensionality Reduction	Principal component analysis (PCA), Chi-square, Entropy, Information gain, t-SNE, Symmetric Uncertainty, Manifold learning
Clustering	K-means, Kohonen's self-organized map, Density-based scan, Hierarchical grouping
Data visualization *	Silhouette plot, scatter plot, heatmap, box plot, clusters, t-SNE

Shapes & dimensionality of datasets

 $|\mathbf{I}|$

- A dataset is $a|A| \times |I|$ matrix, where a|A| is the number of Attributes and |I| is the number of instances
 - Vertical datasets are |A|<<|I|
 - Horizontal datasets are |A|>>|I|
 - Usually, classification methods require a large amount of instances for training a model
 - A large number of atributes can difficult to find a suitable model, mainly those for descriptive analysis



What is the problem ?

- How to choose the best atributes of a given dataset for constructing a predictive or descriptive model ?
- If the number of attributes is small:
 - If the user has previous knowledge about the data domain, just select the more representative atributes
 - If the computational power is large, perform a factorial experimente, that is, test all possible combinations of attributes and compare results.
- If the number of attributes is (very) large:
 - It is imperative to use a computationally efficient method to select the most adequate attributes and discard the remaining



Further reasons for doing attribute selection

- Real-world datasets can have:
 - Irrelevant attributes: those which are not correlated with the target-atribute (class)
 - Noisy attibutes: those with variated degrees of noise due to the collection procedure
 - **Redundant attributes** : those strongly correlated with other atributes
- Using all attibutes can lead to poor predictive models
- Using less attibutes can:
 - Facilitate the visualization of multidimensional data
 - Create models more comprehensible for the user
 - Decrease training time of the algorithms
 - Increase the generalization capability of the training models

Feature selection X Dimensionality Reduction

- Both aim at decreasing the number of attibutes (=features) c dataset so that machine learning methods can be more efficient
- Attribute Selection:
 - Methods that select a **subset** of the original atributes by means of a specific quality criterion
- Dimensionality Reduction:
 - Methods that **create** other atributes by using the original ones



Attribute Selection modes

- Defining a dataset (DS) as a collection of n instances and matributes ($n \times m$ matrix) an atribute selection method simply creates a mapping Φ =DS \rightarrow DS' such that m' << m
- Attibute selection can be performed in three diferente ways:
 - Selecting a subset of the original atributes, within ALL possible combinations (factorial experiment)
 - Ordering the atributes in decreasing order of a measure of suitability to the problem, and establishing a cuttof point
 - Weighting the atributes according to their suitability for a specific task (e.g. classification)

Groups of Methods for Feature Selection

• <u>Filter</u> methods:

- Use statistical measures to measure a "quality" value for each atribute
- Attributes are ordered according to this "quality" criterion
- These methods are most commonly used for **unsupervised** problems

• <u>Wrapper</u> methods:

- Use a search method do swap the possible space of attribute combinations
- The quality of each set of selected attributes is evaluated by a predictive model (e.g. a classifier)
- These methods are computationally expensive
- Most commonly used for supervised problems

Filter methods for feature selection

- Information Theory-based methods:
 - Entropy (H)
 - Information Gain (IG)
 - Average Symmetric Uncertainty (ASU)
 - Minimum-redundance-Maximum-relevance (mRMR)
 - Relief
 - Fast Filter for Feature Selection (FCBF)
- Statistics-based methods:
 - Pearson correlation
 - Spearman correlation
 - Chi-Square Score (Mann-Whitley)
 - Z-Score (Signal-to-Noise Ratio)
 - Between and Within-Group Sum of Squares Ratio (BWSS)
 - Kruskall-Wallis (KW) score

Information Theory-based methods

• Entropy (H): it is a measure of uncertainty, dispersion or disorder, and varies between $0 \le H(x) \le 1$

$$H(X) = -\sum_{x \in X} p(x) \log(p(x))$$

 Information Gain (IG), Mutual Information (MI) or Kullback-Leibler Divergence (KLD): it is the information gain due to the choice of a specific attribute (or set of) with respect to the class

$$IG(X,Y) = H(X) + H(Y) - H(X,Y)$$

• Symmetric Mean Uncertainty (SU): it is derived from IG and corrects possible distortions of this method $SU(X,Y) = 2 \frac{H(X) + H(Y) - H(X,Y)}{H(X) + H(Y)}$

Fast Correlation-Based Filter (FCBF) algorithm

- FCBF selects atributes that are highly correlated with the target atribute (class) and have little (or no) correlation with the other variables. Steps:
 - The correlation measure used is the <u>Symmetric Uncertainty (</u>SU)
 - FCBF selects atributes correlated with the class above a given threshold
 - Then it detects the <u>predominant correlations</u>: when the correlation of an atribute A with the class is greater than the correlation of any other atribute with A
- FCBF was originally proposed for discrete attributes. For continuous attributes a discretization procedure must be applied prior to the use of the algorithm
- FCBF is efficient for high-dimensionality datasets

Yu, L. & Liu, H., Feature Selection for High-Dimensional Data: A Fast Correlation-Based Filter Solution. Proc. 20th Int. Conf. Machine Learning, 2003 https://jundongl.github.io/scikit-feature

Relief algorithm

- It is a "family" of filter methods for feature selection, inspired in Instance-Based Learning (IBL).
- Variants: Relief, ReliefF, RReliefF, I-Relief, tuned-Relief
- The algorithm is based on a user-defined relevance vector *R* and a threshold *L*
- Relief evaluates the differences between pairs of closest neighboring instances (in *n*-dimensional space):
 - If instances of the same class are close (regarding the threshold *L*) the *R* score decreases
 - If they are of diferente classes, the *R* score increases
- Each element of vector *R* corresponds to the relevance of an atribute for the class

Kira, K., Rendell, L.A. The feature selection problem: traditional methods and a new algorithm. Proc. AAAI conference, 1992

https://medium.com/@yashdagli98/feature-selection-using-relief-algorithms-with-python-example-3c2006e18f83

Comparison between FCBF and Relief

	FCBF	Relief
Metric	Correlation	Neighborhood
Computational complexity	O (instances ² x atributes)	O (instances ² x atributes x iterations)
Data type	Discrete or discretized data	Nominal or numeric
Advantages	 Good for high-dimensional data Considers both redundance and relevance 	 Efficient when there are iteractions between atributes Robust to noise
Disadvantages	 Based on linear correlation, may not capture nonlinear correlations 	 Do not remove redundant atributes Performance dependent on the choice of #neighbors

Wrapper methods for feature selection

- They are computationally expensive because they sequentially evaluate a large number of attribute subsets.
- Search models:
 - Random Search, Exhaustive Search (Grid Search),
 - Genetic algorithm,
 - Swarm methods: PSO, ACO, ABC..
 - Classical metaheuristics:Tabu Search, simulated annealing, scatter Search
- Classification methods:
 - OneR, Decision trees (C4.5), Neural network, SVM, KNN
- Quality metrics with cross-validation:
 - Accuracy (for dataset with balanced classes), F1, AUC

Case study #1: Pima indians diabetes – Filter methods

0.13884

0.1398

- Instances: 768 (500 neg., 268 pos.)
- Attibutes: 8 (Number of pregnances Plasma glucose concentration, Dias blood pressure, Skin tickness, 3-hou serum insulin, Body-mass index, Dia pedigree function, Age)

4 SkinThickness

Classes: Positive, Negative

					#	Info	. gain	Gain ratio	χ²	Re	liefF	FCBF 💌		
s pos.)		1 N G	ilucose			_	0.170	0.085 139.90		1	0.019	0.131		
nances,		2 N A	ge				0.081	0.041	62.02	9.	0.005	0.059		
n, Diastolio	2	в Ŋ В	MI			0.079		0.039	53.74	4	0.010	0.057		
s, 3-nour day Diaba	tos	4 N C	DiabetesPedigreeFur	nction			0.022	0.011	16.14	3	0.005	0.015		
uex, Diabeles		5 Ŋ Ir	nsulin		0.055		0.030	8.78) 	0.009	0.000			
6 N Pregnancies						_	0.043	0.021	0.021 34.316		0.014	0.000		
		7 N S	kinThickness			-	0.036	0.018	5.262		0.010	. 0.000		
	-	в Ŋ В	loodPressure			-	0.015	0.007	12.91	в	0.012	0.000		
						•		•	•					
Feature	Info. g	jain	Gain ratio		χ²	R		ReliefF	FCBF		a	verage		
Glucose		1.000	1.00			1		1.00		1.000	1.000			
BMI		0.41088	0.4108		0.36	0092		0.4562	0	.43144		0.341876		
Aae		0.42774	0.4284		0.42	1625	_	0.0783	0	.44712		0.276309		
Pregnancies	_	0.17826	0.1807	_	0.21	0.215796		0.7983	0	.00002	0.231446			
BloodPressure		0.000	0.00	-	0.056	8603	_	0.8275		0.000	_	0.165493		
Insulin	_	0.25635	0.2926		0.026	1313		0.0391	0.00002		-	0.117607		
DiabetesPedigr		0.04639	0.0463	_	0.080	3189		0.3687	0.11583		_	0.11544		

0



**Average of Info.Gain, Gain Ratio, ReliefF, FCBF

0.00

0.00001

0.0557263

Case study #1: Pima indians diabetes – Filter methods

- Classification using decision tree and neural network
- 10-fold cross-validation
- Evaluation metric: F1 (unbalanced classes) and tree complexity

Features	Nodes/leaves	Classification Results
Top-1	103 / 52	Model AUC CA F1 Prec Recall MCC
(OneRule)		Tree (1) 0.741 0.715 0.698 0.703 0.715 0.332
		Neural Network 0.785 0.745 0.731 0.738 0.745 0.408
Top-3	91 / 46	Model AUC CA F1 Prec Recall MCC
		Tree (2) 0.793 0.762 0.759 0.758 0.762 0.465
		Neural Network 0.829 0.768 0.762 0.763 0.768 0.473
All features	91 / 46	Model AUC CA F1 Prec Recall MCC
		Tree 0.748 0.733 0.730 0.728 0.733 0.400
		Neural Network 0.824 0.758 0.755 0.754 0.758 0.458

All atributes \rightarrow worse re

Case study #1: Pima indians diabetes – Wrapper methods

- Orange does **not** have Wrapper methods for feature selection !
- Classifier + Search method (Weka)



20



Case study #2: Brain tumor images classificaton

- Magnetic resonance images of three classes of tumors (glioma, meningioma, pituitary) plus a control group (no-tumor)
- Data available at: <u>https://www.kaggle.com/datasets/sartajbhuvaji/brain-</u> <u>tumor-classification-mri</u>
- Feature extractor: Squeezenet NN (1000 numeric features)
- Train / test: 2870/364 images





Case study #2: Brain tumor images classificaton

• All runs use 5x stratified cross-validation in the training dataset

								-							
All1000	Model	AUC	CA	F1	Prec	Recall	MCC	Top-1	Model	AUC	CA	F1	Prec	Recall	MCC
features	Tree Baseline-All	0.771	0.686	0.685	0.684	0.686	0.571	Informatio	Tree Baseline Top-1	IG 0.59	2 0.362	0.362	0.369	0.362	0.127
	Neural Network	0.983	0.903	0.903	0.903	0.903	0.867	n Gain	Neural Network	0.67	5 0.434	0.402	0.442	0.434	0.216
Top-23	Model	AUC	CA	F1	Prec	Recall	MCC	Top200	Model	AUC	CA	F1	Prec	Recall	мсс
FCBF	Tree-23 FCBF	0.785	0.684	0.685	0.687	0.684	0.568	Informatio	Tree Top-200-IG	0.781	0.693	0.692	0.691	0.693	0.580
	Neural Network	0.960	0.831	0.831	0.831	0.831	0.769	n Gain	Neural Network	0.979	0.886	0.886	0.886	0.886	0.845
Top-100	Model	AUC	CA	F1	Prec	Recall	MCC	Top200	Model	AUC	CA	F1	Prec	Recall	MCC
Relief	Tree Top-100 Relief	0.777	0.68	5 0.68	5 0.685	5 0.685	0.570	Informatio	Tree Top-200 IGR	0.780	0.692	0.691	0.690	0.692	0.579
	Neural Network	0.975	0.87	3 0.87	3 0.873	3 0.873	0.826	n Gain	Neural Network	0.979	0.892	0.892	0.892	0.892	0.853
								Ratio							

Dimensionality Reduction

- It is **not** an attribute selection method, but a transformation of the original attributes into new and, hopefully, more discrimatory ones
- Main methods:
 - Principal Component Analysis (PCA)
 - Singular Vector Decomposition (SVD)
 - Independent Component Analysis (ICA)
 - Locally Linear Embedding (LLE)
 - Isometric Feature Mapping (ISOMAP)
 - Non-negative Matrix Factorization (NMF)
 - Autoencoders
 - Some of these methods are available in the Pyhton's scikit-learn library:

https://machinelearningmastery.com/dimensionality-reduction-algorithms-with-python/



- Also known as Karhunen-Loève transform, but invented by Pearson (1901)
- It is an orthogonal linear transformation method that aims to find new linearly uncorrelated atributes.
- Besides Dimensionality Reduction, the method is also used for Exploratory Data Analysis tasks
- Considering *D* as the number of original atributes, PCA's computational complexity is given by:
 - Time: **O**(D³),
 - Memory: **O**(D²)

- Each Principal Component (PC) represents a small amount of the **variance** in the data
- The first PC captures the largest possible variance in the data
- Each subsequent component has the maximum remaining variance *, considering the constraint of being orthogonal (uncorrelated) to all the previous ones
- Usually, a number of componentes is chosen such that 90-95% of the original variance is explained (when possible)

* Variance: it is a measure of statistical dispersion, representing how far its values are from the expected value

• Variance of a single variable:

$$\mathrm{Var}(X) = rac{1}{n}\sum_j (ar{x}-x_j)^2 = \sigma_X^2$$

• Covariance between two variables: Cov

$$\mathrm{Cov}(X,Y) = rac{1}{n}\sum_j (ar{x}-x_j)(ar{y}-y_j) = \sigma_{XY}$$

• Covariance matrix for **n** variables:

$$\mathrm{C} = egin{pmatrix} \sigma_{11} & \sigma_{12} & \ldots & \sigma_{1n} \ \sigma_{21} & \sigma_{22}^2 & \ldots & \sigma_{2n} \ dots & dots & \ddots & dots \ \sigma_{n1} & \sigma_{n2} & \ldots & \sigma_{nn}^2 \end{pmatrix}$$

https://wilkelab.org/SDS375/slides/dimension-reduction-1.html

- Diagonalization of the covariance matrix: $\mathbf{C} = \mathbf{U}\mathbf{D}\mathbf{U}^{T}$
- Where:
 - U = rotation matrix, D = diagonal matrix
 - λ_j^2 = the variance explained by the *j* component (or eigenvalues)
 - The componentes are uncorrelated
 - The covariance between the componentes is zero

$$\mathbf{U} \begin{pmatrix} \lambda_{1}^{2} & 0 & \dots & 0 \\ 0 & \lambda_{2}^{2} & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \lambda_{n}^{2} \end{pmatrix} \mathbf{U}^{\mathrm{T}}$$

- The first componente (PC1) accounts for the dimension where there is the greatest variance in the data
- The greater the number of PCA compoenents, the greater the **cumulative variance** "explained"



Case study #3: Bluejay dataset

- IMPORTANT: this is a **trivial** problem, since there are only 6 variables
- Variables (6): BillDepth, BillWidth, BillLenght, Head, Mass, Skull
- Meta-attribute: Sex (0=female, 1=male)





Case study #3: Bluejay dataset

• Numer of instances: 124

bird_id	sex	bill_depth_mm	bill_width_mm	bill_length_mm	head_length_mm	body_mass_g	skull_size_mm
0000-00000	м	8.26	9.21	25.92	56.58	73.30	30.66
1142-05901	М	8.54	8.76	24.99	56.36	75.10	31.38
1142-05905	М	8.39	8.78	26.07	57.32	70.25	31.25
1142-05907	F	7.78	9.30	23.48	53.77	65.50	30.29
1142-05909	М	8.71	9.84	25.47	57.32	74.90	31.85
1142-05911	F	7.28	9.30	22.25	52.25	63.90	30.00
1142-05912	М	8.74	9.28	25.35	57.12	75.10	31.77
1142-05914	м	8.72	9.94	30.00	60.67	78.10	30.67
1142-05917	F	8.20	9.01	22.78	52.83	64.00	30.05
1142-05920	F	7.67	9.31	24.61	54.94	67.33	30.33
1142-05930	М	8.78	8.83	25.72	56.54	76.40	30.82

Case study #3: Bluejay dataset

- PCA aligns the major axes with directions of maximum variation in the data
- Using PCA, the 6-dimension original space is represented by its principal components with the highest variance



Case study #4: Smoking effects on B lymphocytes dataset

- The B lymphocytes attack invaders outside the cells
- They are responsible for adaptive immunity and are associated with the onset and development of many diseases induced by continuous cigarette use.
- The analysis of gene expression of these cells may provide useful information about the relationship between the cells and diseases



Case study #4: Smoking effects on B lymphocytes dataset

- Gene expression data from peripheral blood B cells of 39 smokers and 40 non-smokers, all female
- The data are from a sample of 3000 genes out of a total of 14,500 genes evaluated

Pan F. et al. Impact of female cigarette smoking on circulating B cells in vivo: the suppressed ICOSLG, TCF3, and VCAM1 gene functional network may inhibit normal cell function. Immunogenetics, 62(4), 237-251, 2010. https://www.ncbi.nlm.nih.gov/sites/GDSbrowser/ ?acc=GDS3713



Case study #4: Smoking effects on B lymphocytes dataset Dimensionality reduction using PCA

Σ

39

40

79



Actual

# PC's	%varian ce	F1 DT	F1 NN
No PC		0,747	0,911
1 PC	8%	0,582	0,620
2 PC	14%	0,822	0,848
10 PC	34%	0,835	0,899
20 PC	49%	0,835	0,873
30 PC	61%	0,848	0,835
50 PC	80%	0,835	0,543
79 PC	99%	0,835	0,187

Case study #4: Smoking effects on B lymphocytes dataset Feature selection #Attrib. F1 DT F1 NN Reduced Data Neural Network (2) Selected Data -> Data £ 213 ¥ = = baseline 0,911 all 0,747 Datasets (1) Select Columns (2) Rank Data Table Test and Score (1) Confusion Matrix Tree (1) Tree Viewer (3)

- Depending upon the classifier algorithm, the number of attributes (variables) can influence positivelly or negativelly in the classification result
- Comparison between: Relief, FCBF, InfoGain, and χ^2

Relief	50	0,810	0,975
Relief	100	0,733	0,987
FCBF	25	0,706	1,000
FCBF	100	0,734	0,949
InfoGain	25/100	0,734	0,911
X^2	25	0,747	0,925
X^2	100	0,732	0,937

Case study #2 (revisited):

Brain tumor images classification

• Using PCA for creating up to 90 new features from the original 1000, but using only 2

Method	(previous best result) Top-200 features (selected by Information Gain)								2 PCs						
Results in the test set	Model Tree Top-200 IGR Neural Network	AUC 0.780 0.979	CA 0.692 0.892	F1 0.691 0.892	Prec 0.690 0.892	Recall 0.692 0.892	MCC 0.579 0.853		Model Tree 30PC Neural Network	AUC 0.771 0.983	CA 0.688 0.907	F1 0.687 0.908	Prec 0.686 0.908	Recall 0.688 0.907	MCC 0.574 0.874
Best ever result !															

Additional links for further study

• Feature selection in Python using Relief for feature selection

https://medium.com/@yashdagli98/feature-selection-using-relief-algorithms-with-pythonexample-3c2006e18f83

- Feature selection in Python using FCBF for feature selection: <u>https://github.com/shiralkarprashant/FCBF</u>
- Feature selection methods using Scikit-learn https://scikit-learn.org/stable/modules/feature_selection.html
- A comprehensive guide on Feature Selection (Kaggle): <u>https://www.kaggle.com/code/prashant111/comprehensive-guide-on-feature-selection</u>