

Data Mining
Final project report
Implementation of a single classification algorithm and several
algorithms for feature reduction.

Przemysław Piórkowski, Michał Debski, Damian Hermanowski

June, 2015

Contents

1	Introduction	2
2	Algorithms implementation	3
3	Experiments	6
4	Conclusions	9

Chapter 1

Introduction

The project objective is to implement single classification algorithm and several algorithms for feature reduction. In the addition, developed algorithms have to be evaluated against chosen data set and classification method.

Basing on the document provided by project's Supervisor and the initial research the Team decided to implement three types of the feature reduction algorithms: Document Frequency Thresholding (DFT), Mutual Information (MI) and χ^2 statistic (CHI). As a classification algorithm the k-nearest neighbors (kNN) method was chosen. The Team decided to use the C++ programming language for implementation purposes.

The source data set: SPAM E-mail Database, consists of 4601 instances of e-mails. Each e-mail is considered either as a spam or not. Classification is based on 57 continuous attributes denoting particular words and characters frequency as well as lengths of uninterrupted sequences of capital letters in a message. The data set does not have any missing attribute values.

Chapter 2

Algorithms implementation

DFT thresholding is the simplest technique for attribute reduction. Document frequency is the number of documents in which a term occurs. Our implementation counts occurrences of not null values of each attribute. It utilizes the training set for this step. After that a vector of frequencies is created. An attribute which frequency is below given threshold is filtered out for further classification steps. The basic assumption is that rare terms are either non-informative for category prediction, or not influential in global performance.

MI tries to measure the influence of existence of particular attribute/term on a given class value. If one considers the twoway contingency table of a term t and a category c , where A is the number of times t and c co-occur, B is the number of time the t occurs without c , C is number of times c occurs without t , and N is the total number of documents, then the mutual information criterion between t and c is estimated using formula:

$$I(t, c) \approx \log \frac{A \times N}{(A + C) \times (A + B)}$$

$I(t, c)$ equal to zero means that attribute/term t and category c are independent. The goodness of aterm/attribute is scored in two ways: as a average of MIs for a given attribute (ML_AVG mode) and as a maximal MI for a given attribute (ML_MAX mode). The outcome of the algorithm is a reduction vector of attributes scores. An attribute with a score below a given threshold is filtered out.

$$I_{avg}(t) = \sum_{i=1}^m P_r(c_i) I(t, c_i)$$

$$I_{max}(t) = \max_{i=1}^m \{I(t, c_i)\}$$

The χ^2 statistic measures the lack of independence between an attribute t and a class c . Similarly to MI algorithm utilizes the twoway contingency table with an additional value D , which equals to case when neither c nor t occurs. The term-goodness is calculated with a given formula:

$$\chi^2(t, c) = \frac{N \times (AD - CB)^2}{(A + C) \times (B + D) \times (A + B) \times (C + D)}$$

Similarly to the MI method, the CHI statistic between the given attribute and classes may be computed as two scores: average value (CHI_AVG mode) and maximum value (CHI_MAX mode). The outcome of the algorithm is a vector of attributes scores. An attribute with a score below a given threshold is filtered out.

$$\chi_{avg}^2(t) = \sum_{i=1}^m P_r(c_i) \chi^2(t, c_i)$$

$$\chi_{max}^2(t) = \max_{i=1}^m \{\chi^2(t, c_i)\}$$

The project implements three feature reduction algorithms within 3 classes:

- CHIReduction
- DFTReduction
- MIReduction

The kNN algorithm is implemented in the **KNNClassifier** class. There were developed additional classes for data loading and data interpretation: **DataLoader**, **Data**. **DataAdapter** class is used to iterate over data in function defined manner. There is the **Logger** class for debugging and echoing purposes. There is the class supporting multithreading for parallelization capabilities: **ParallelExecutor** with a support of *AtomicIfPossible* and *AtomicInternal* classes.

The developed project application can be run with given parameters:

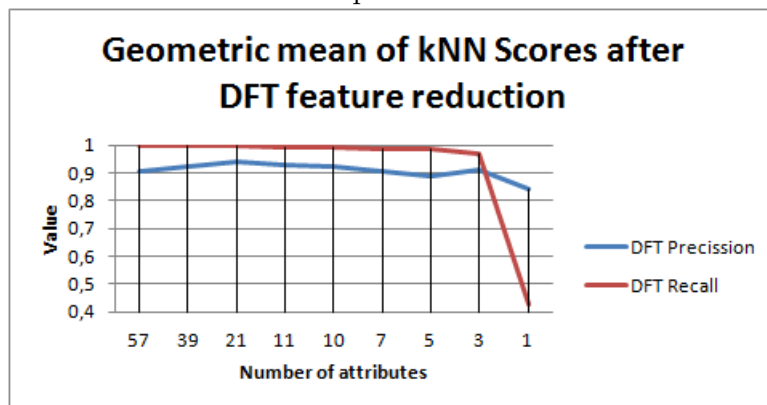
- `-ftrain trainFileName`
- `-ftest testFileName`

- `-h` - switch for reading headers of a source file
- `-cname className`
- `-k value`
- `-CV value` - number of iterations/folds of cross validation
- `-DFT dftThreshold`
- `-MI_MAX` - switches on MAX score calculation for MI algorithm
- `-MI_AVG` - switches on AVG score calculation for MI algorithm
- `-MI miThreshold`
- `-l verbosityLevel`
- `-CHI_MAX` - switches on MAX score calculation for CHI algorithm
- `-CHI_AVG` - switches on AVG score calculation for CHI algorithm
- `-CHI chiThreshold`

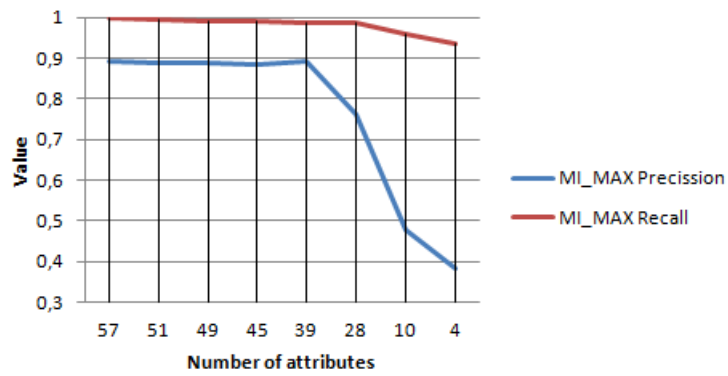
Chapter 3

Experiments

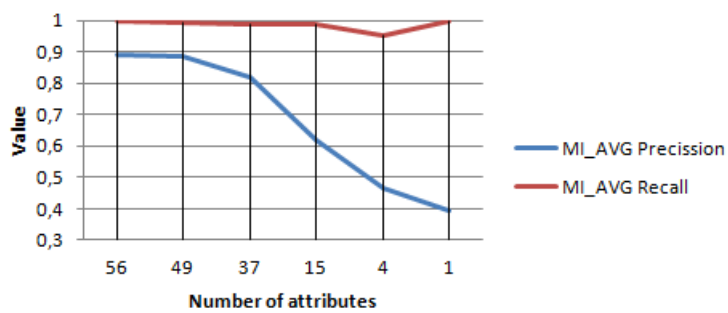
In order to validate implemented algorithms against the data set the *7-fold Cross-validation* method was used (-CV 7). kNN implementation has one control parameter - the number of neighbors voting. For evaluation purposes this k-value was set to 5. The influence of each feature reduction algorithm on the kNN classification method was measured with precision and recall values. Scores were calculated as a geometric mean of values obtained from each round of the cross validation. These quality measures were compared with the results of classification without feature reduction. Number of attributes was controlled with picked threshold values.



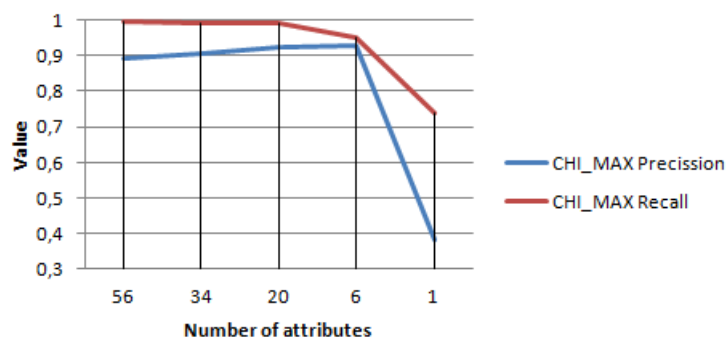
**Geometric mean of kNN Scores after
MI_MAX feature reduction**



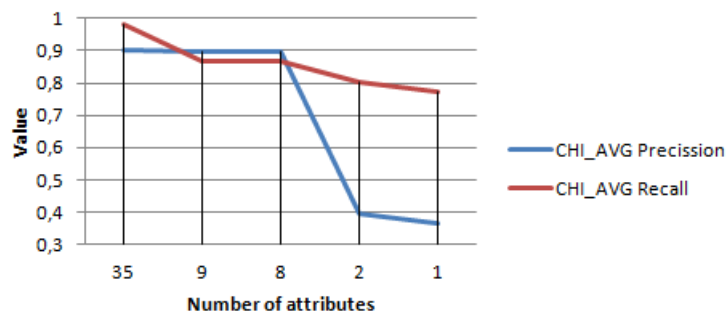
**Geometric mean of kNN Scores after
MI_AVG feature reduction**



**Geometric mean of kNN Scores after
CHI_MAX feature reduction**



**Geometric mean of kNN Scores after
CHI_AVG feature reduction**



Chapter 4

Conclusions

The results are interesting. The application of DFT, CHI_MAX and MI_MAX on the given data set gave slight precision gain when less than a half of attributes were cut. The DFT method achieved high precision and recall scores even with 3 attributes left - it is a great result for a such simple algorithm. CHI_MAX and CHI_AVG methods marked comparable - good overall results. Mutual Information (MI) methods showed worse results for aggressive feature reduction.