

Classification of Email Spam using Concept Learning(Group 53)

VAMSI SRI NAGA MANIKANTA MURUKONDA
20010620-T114
vamu21@student.bth.se

SREE LAKSHMI HIRANMAYEE KADALI
20010920-T244
srkd21@student.bth.se

I. INTRODUCTION

Based on concept learning, algorithms 4.1 and 4.3 from the main literature is implemented for this project work. We use both the algorithms to create a conjunctive least general generalization algorithm (LGG) to classify spam base dataset to distinguish between spam and non-spam mails by returning the most general version space.

II. PREPROCESSING OF THE DATA

The dataset given consists of 4601 instances with 57 attributes, with the last attribute being the grouping of mail as spam(1) and non-spam(0). The dataset is classified into two sets "train_data" and "test_data" using the sklearn and pandas library. The instances are read into one dataset "train_data" and classification values (1-spam, 0-non spam) is read into other dataset "test_data". The datasets are then transformed into numpy arrays which are then transformed to train and test sets with the ratio of 70:30 where, 70% of the spam mails are used to train and 30% is used to test. Since, the data is discretized separately for training and testing, fewer bins are used to provide greater accuracy, and n=10 bins are used. Thus KbinsDiscretizer is used. Later the datasets are mixed with non-spam emails. A hypothesis space is then created and data is calculated into the hypothesis space.

III. ALGORITHM SELECTION AND IMPLEMENTATION

A combination of algorithms 4.1 and 4.3 is used to give better accuracy for the classification of data. The LGG algorithm is employed as soon as each occurrence in the feature set is set up in the hypothesis space. Each one of the algorithms, LGG's internal disjunction and LGG's conjunction receive occurrences one by one and examine the hypothesis space for its particular attribute set. If the hypothesis attribute set, as of now doesn't contain the attribute, the attribute is added and the iteration is repeated until the last instance. This modified hypothesis space return the required hypothesis space.

IV. TRAINING

For the training purpose, 70% of the dataset with spam mails is supplied to the LGG algorithm with LGG conjugate

inner disjunction. The implementation builds a parameter space which contains lists of internal disjunction sets for each feature.

V. TESTING

The test dataset which contains 30% of the spam messages and rest as non spam messages is provided for the algorithm for testing. Each occurrence is assessed at a time to confirm it is inside the hypothesis space. The occasion set is classified according with the hypothesis space as spam and non-spam.

VI. HYPOTHESIS SPACE

The calculated hypothesis space as defined in the main literature section 4.1 is with the number of possible extensions i.e set of instances $2*(10*57)$. The possible instances with a conjugate concept includes $(10+1)*57$ including the absence of feature as an addition 'value'.

[illegible]

VII. RESULTS

The accuracy of the model is given by $(tp+tn)/(tp+tn+fn+fp)$
The accuracy of the model is 0.52.