

Evaluation of email classification methods

The problem statements

This study is about designing a spam filter that can separate ham and spam emails based on various machine learning techniques. The subject line of the emails is analyzed based on the NLP technique of tokenizing sentences. Two methods of Machine learning algorithms are used (1) dimensionality reduction technique PCA is applied to Logistic regression and (2) 1D convolution neural network to train a model on training set and then asses it on a test set. Finally, the models are evaluated on bigdata set of 58000 emails and accuracy and precision is tested.

Dataset

The dataset is based on cleaned Enron corpus, there are a total of 92188 messages belonging to 158 users with an average of 757 messages per user. The dataset has almost an equal distribution of ham and spam emails. In this study 2000 emails are used 1000 ham + 1000 spams. Each email text is preprocessed with python library mail parser to extract various features. The features extracted are shown in table 1.

Parameter	Example
Body	"Carolina Power & Light and Florida Power Corporation submitted tariff revisions in compliance with a Commission order addressing the energy imbalance provisions that apply in CP&L's zone. \n\nThe proposed revisions include:\n\t-return-in-kind provisions for energy imbalances\n\t- deletion of the separate capacity charge for undersupply of energy outside the deadband for 10+ hours in a month\n\t- a provision that deficient energy will be offset or credited with energy associated with spinning and supplemental reserves.\n\nInterventions/protests are due Aug. 15.\n\nIf you would like further information please contact me.\n\nSusan Scott Lindberg\n\n30596"
Date	2001-08-05 18:34:50
From	[(Scott, Susan, Susan.Scott@ENRON.com)]
To	[('Acevedo, Rudy', 'Rudy.Acevedo@ENRON.com'), ('Carson, Mike', 'Mike.Carson@ENRON.com'), ('Comeaux, Keith', 'Keith.Comeaux@ENRON.com'), ('Connor, Joe', 'Joe.Connor@ENRON.com'), ('Fairley, David', 'David.Fairley@ENRON.com'),]
message_id	'<902B8E00B151D44C98CA48BDD4BEA3F5082C01@NAHOU-MSMBX01V.corp.enron.com>'
subject	'CP&L tariff changes (ER01-1807)'

Table 1. Parameters extracted from emails. Subject line is used in this study

Processing on email subjects

The preprocessing of email subject lines includes tokenize the email subject lines. The python library spacy is used to in this case. This method counts the number of times for the occurrences of tokens. The dataframe thus created has > 3500 columns based on the individual tokens identified by the algorithm. In this problem there are more attributes than the number of rows. However, since the dimensionality of the problem is very high its needs to be reduced to lower dimensions by *latent semantic indexing*.

PCA with logistic regression

PCA reduces the dimensionality hence the complexity while maintaining structure (variance) of a dataset. It performs a rotation of the data that maximizes the variance in the new axes.

By combining the advantages of both the PCA and logistic regression, PCA is used to extract feature and reduce the dimensions of process data. Afterwards logistic regression is used as the classifier for spam and ham emails.

The scatter matrix is used to extract feature, and then obtain all the individual characteristics subspace $W_i, i = 1, 2, \dots, m$. First the eigenvalues and the corresponding eigenvectors is used to generate matrix. Second the eigenvalues are order from largest to smallest, and similarly putting the corresponding eigenvectors in order of largest

to smallest, the optimal projection matrix (X_1, X_2, \dots, X_d) is thus created and is

associated with the d largest generalized eigenvalues. Finally, logistic regression will be used as the classifier of ham and spam emails.

Logistic regression is used classify the emails based on > 3500 attributes. A set of regularization parameter in C is used to reduce overfitting, $C=0.1$ performs the best on both train-test set.

PCA and logistic regression combination are tested with 3 different samples sizes 16, 160 and 1600. Sample size of 160 have best accuracy and precision values on training and test set.

Method	Parameters	Training Set						Test Set					
		True_pos	False_pos	False_neg	True_neg	Accur	Preci	True_pos	False_pos	False_neg	True_neg	Accur	Preci
logistic regression	C=0.01	519	19	281	781	0.81	0.65	109	11	91	189	0.75	0.55
logistic regression	C=0.1	608	15	192	785	0.87	0.76	126	23	74	177	0.76	0.63
logistic regression	C=1	759	11	41	789	0.97	0.95	146	32	54	168	0.79	0.73
logistic regression	C=10	800	17	0	783	0.99	1.00	182	91	18	109	0.73	0.91
logistic regression	C=100	800	16	0	784	0.99	1.00	183	100	17	100	0.71	0.92
PCA+ logistic regression	Feature_column = 16	480	37	320	763	0.78	0.60	111	11	89	189	0.75	0.56
PCA+ logistic regression	Feature_column = 160	629	45	171	755	0.87	0.79	140	34	60	166	0.77	0.70
PCA+ logistic regression	Feature_column = 1600	800	17	0	783	0.99	1.00	181	91	19	109	0.73	0.91