# Machine learning basics with applications to email spam detection

**Lydia Song, Lauren Steimle, Xiaoxiao Xu, and Dr. Arye Nehorai**
Department of Electrical and Systems Engineering, Washington University in St. Louis

Washington University in St.Louis — SCHOOL OF ENGINEERING & APPLIED SCIENCE

## Abstract

Machine learning is a branch of artificial intelligence concerned with the creation and study of systems that can learn from data. A machine learning system could be trained to distinguish between spam and non-spam (ham) emails. We aim to study current methods in machine learning to identify the best techniques to use in spam filtering. We found that the One-Nearest Neighbor algorithm achieved the best performance.

## Motivation

Email has become one of the most important forms of communication. In 2013, there were about 180 billion emails sent per day worldwide and 65% of the emails sent were spam emails. Links in spam emails may lead to users to websites with malware or phishing schemes. Therefore, an effective spam filtering technology is a significant contribution to the sustainability of the cyberspace and to our society.

Many spam filters rely on Domain Name System-Based Blackhole Lists to keep track of IP addresses that send large amounts of spam so that future email from these addresses can be rejected. However, spammers are circumventing these lists by using larger numbers of IP addresses. Current blacklisting techniques could be paired with content-based spam filtering methods to increase effectiveness.

## Methods

Machine learning systems operate in two stages: training and classification.



Image credit: "Statistical pattern recognition: A review"

The goal of preprocessing is to remove any noise and normalize the data, and create a compact representation of the data. In training mode, the classifier will determine input patterns from a set of training data and determine how to partition the feature space. In testing mode, the classifier assigns testing data to a class based on their features. Performance results are determined from these classifications.

A dataset of 1000 emails from the Text Retrieval Conference (TREC) 2007 corpus was used to train and test the classifiers.

## Preprocessing

**Spam Email in Web Browser**

Your history shows that your last order is ready for refilling.

Thank you,

Sam Mcfarland
Customer Services

**Spam Email in Data Set**



**4. Create a feature matrix**

|        | 'histori' | 'last' | ... | 'servi' |
|--------|-----------|--------|-----|---------|
| Email 1 | 1 | 1 | ... | 1 |
| Email 2 | 0 | 3 | ... | 5 |
| ⋮ | ⋮ | ⋮ | ⋱ | ⋮ |
| Email m | 2 | 1 | ... | 6 |

**1. Tokenize**

tokens= ['your', 'history', 'shows', 'that', 'your', 'last', 'order', 'is', 'ready', 'for', 'refilling', 'thank', 'you', 'sam', 'mcfarland', 'customer services']
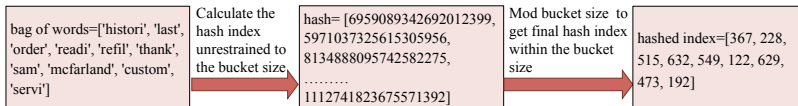
**2. Filter words**

filtered_words=['history' 'last', 'order', 'ready', 'refilling', 'thank', 'sam', 'mcfarland', 'customer', 'services']

**3. Bag of words**

bag of words=['histori', 'last', 'order', 'readi', 'refil', 'thank', 'sam', 'mcfarland', 'custom', 'servi']

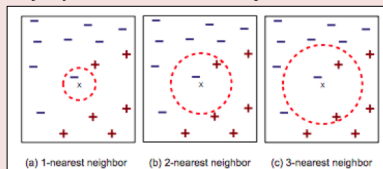## Dimensionality Reduction - Hashing Trick

Without Hashing, the dimensionality of the feature matrix with 70 emails is 9403. After Hashing, the dimensionality is reduced to the number of hash buckets (572, 1024, or 2048).
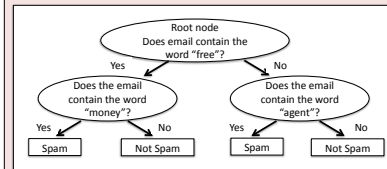
bag of words=['histori', 'last', 'order', 'readi', 'refil', 'thank', 'sam', 'mcfarland', 'custom', 'servi']

→ Calculate the hash index unrestrained to the bucket size →

hash= [6959089342692012399, 5971037325615305956, 8134888095742582275, ………… 1112741823675571392]

→ Mod bucket size to get final hash index within the bucket size →

hashed index=[367, 228, 515, 632, 549, 122, 629, 473, 192]

## Classifiers

### k-Nearest Neighbors

Predict the label of a data point X by using a majority vote of the k closest data points to X



(a) 1-nearest neighbor  (b) 2-nearest neighbor  (c) 3-nearest neighbor

Source: Seyda Ertekin, MIT Opencourseware, 15.097 Spring 2012. Credit: Seyda Ertekin.

### Decision Tree

Map observations about a data point's features to determine its class. The tree is constructed to maximize the information gain of its decisions.



### Naïve Bayesian

Compare $p(\text{Spam}|F_1,\ldots,F_n)$ and $p(\text{Ham}|F_1,\ldots,F_n)$

$$p(C|F_1,\ldots,F_n)=\frac{p(C)p(F_1,\ldots,F_n|C)}{p(F_1,\ldots,F_n)} \qquad p(C|F_1,\ldots,F_n) \propto p(C)p(F_1,\ldots,F_n|C)$$

Assume that the features are independent of each other

$$\arg\max_c \; p(C=c)\prod_{i=1}^n p(F_i=f_i|C=c)$$

### Logistic Regression

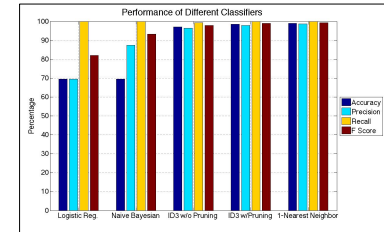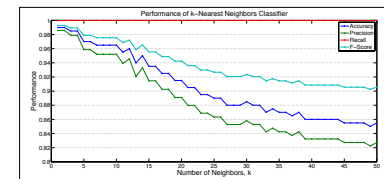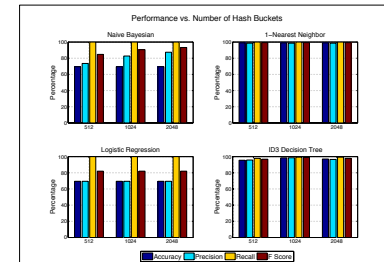Fit a linear model to the feature space $z=\theta_0+\sum_{i=1}^n \theta_i X_i$

$$p(\text{Spam}|z)=\frac{1}{1+e^{-z}}$$

$$p(\text{Ham}|z)=1-p(\text{Spam}|z)=\frac{e^{-z}}{1+e^{-z}}$$

## Results

Key performance measures
• Accuracy – percentage of correctly identified emails
• Precision – percentage of emails classified as spam that were actually spam
• Recall – percentage of spam emails that were accurately classified
• F-score – 2*Precision*Recall / (Precision + Recall)



Performance vs. Number of Hash Buckets



Performance of k-Nearest Neighbors Classifier



Performance of Different Classifiers

## Conclusions

The One-Nearest Neighbor algorithm had the best performance with 99.00% accuracy, 98.58% precision, and 100% recall.

All of the algorithms had very high recall performance and lower precision. This suggests it is easy to classify spam emails correctly and more difficult to classify ham emails correctly.

## Literature Cited

A. Jain, R. Duin, and J. Mao. Statistical pattern recognition: A review. *IEEE Trans*. PAMI, 22(1):4–37, 2000.
A. Ramachandran, D. Dagon, and N. Feamster, "Can DNS-based blacklists keep up with bots?," in CEAS 2006 The Second Conference on Email and Anti-Spam, 2006.
G. Cormack, "Email spam filtering: A systematic review," Foundations and Trends in Information Retrieval, vol. 1, no. 4, pp. 335–455, 2008.
Christina V, Karpagavalli S and Suganya G (2010), 'A Study on Email Spam Filtering Techniques', International Journal of Computer Applications (0975 – 8887), Vol. 12- No. 1, pp. 7-9

## Acknowledgements

## Further Information

We encourage those that are interested in learning more about statistical pattern recognition and machine learning to look to Andrew Ng's Machine Learning course at www.coursera.org.