



Center for Machine Learning and Intelligent Systems

[About](#) [Citation Policy](#) [Donate a Data Set](#) [Contact](#)

Repository Web

[View ALL Data Sets](#)

Spambase Data Set

Download: [Data Folder](#), [Data Set Description](#)

Abstract: Classifying Email as Spam or Non-Spam

From	Subject
CalcedPro...	Get the car of your dreams with CalcedProder Help!
TotalSpam...	How Old Are You Really? - Take the SwagAge Test!
@ Donoffs L...	2 Quick way to make it grow[!]
Beryl membe...	visa p-o-8-4-07
Wardroth...	Special To TheGames Member Offer
Accept Credit...	Process Credit Cards for Zero Up Front Cost
James	Your Pharmacy is!
Quick Cash A...	Get A \$500 Cash Advance
Leland Denny	Brushed w/brushes!
eddye kend	Office IP - [60]
Comp Dept	Get a complimentary Starbucks Gift Card on us
Guadalupe N...	Pay NO attention to the Man Behind the Curtain!
Sunbelt Media	Get ready for monday OCTOBER 20TH

Data Set Characteristics:	Multivariate	Number of Instances:	4601	Area:	Computer
Attribute Characteristics:	Integer, Real	Number of Attributes:	57	Date Donated	1999-07-01
Associated Tasks:	Classification	Missing Values?	Yes	Number of Web Hits:	76020

Source:

Creators:

Mark Hopkins, Erik Reeber, George Forman, Jaap Suermondt
Hewlett-Packard Labs, 1501 Page Mill Rd., Palo Alto, CA94304

Donor:

George Forman (gforman at nospam hpl.hp.com) 650-857-7835

Data Set Information:

The "spam" concept is diverse: advertisements for products/web sites, make money fast schemes, chain letters, pornography..

Our collection of spam e-mails came from our postmaster and individuals who had filed spam. Our collection of non-spam e-mails came from filed work and personal e-mails, and hence the word 'george' and the area code '650' are indicators of non-spam. These are useful when constructing a personalized spam filter. One would either have to blind such non-spam indicators or get a very wide collection of non-spam to generate a general purpose spam filter.

For background on spam:

Cranor, Lorrie F., LaMacchia, Brian A. Spam!
Communications of the ACM, 41(8):74-83, 1998.

- (a) Hewlett-Packard Internal-only Technical Report. External forthcoming.
- (b) Determine whether a given email is spam or not.
- (c) ~7% misclassification error. False positives (marking good mail as spam) are very undesirable. If we insist on zero false positives in the training/testing set, 20-25% of the spam passed through the filter.

Attribute Information:

The last column of 'spambase.data' denotes whether the e-mail was considered spam (1) or not (0), i.e. unsolicited commercial e-mail. Most of the attributes indicate whether a particular word or character was frequently occurring in the e-mail. The run-length attributes (55-57) measure the length of sequences of consecutive capital letters. For the statistical measures of each attribute, see the end of this file. Here are the definitions of the attributes:

48 continuous real [0,100] attributes of type word_freq_WORD
= percentage of words in the e-mail that match WORD, i.e. $100 * (\text{number of times the WORD appears in the e-mail}) / \text{total number of words in e-mail}$. A "word" in this case is any string of alphanumeric characters bounded by non-alphanumeric characters or end-of-string.

6 continuous real [0,100] attributes of type char_freq_CHAR
= percentage of characters in the e-mail that match CHAR, i.e. $100 * (\text{number of CHAR occurrences}) / \text{total characters in e-mail}$

1 continuous real [1,...] attribute of type capital_run_length_average
= average length of uninterrupted sequences of capital letters

1 continuous integer [1,...] attribute of type capital_run_length_longest
= length of longest uninterrupted sequence of capital letters

1 continuous integer [1,...] attribute of type capital_run_length_total
= sum of length of uninterrupted sequences of capital letters
= total number of capital letters in the e-mail

1 nominal {0,1} class attribute of type spam
= denotes whether the e-mail was considered spam (1) or not (0), i.e. unsolicited commercial e-mail.

Relevant Papers:

N/A

Papers That Cite This Data Set¹:



Don R. Hush and Clint Scovel and Ingo Steinwart. Los Alamos National Laboratory Stability of Unstable Learning Algorithms. Modeling, Algorithms and Informatics Group, CCS-3. 2003. [[View Context](#)].

Yongmei Wang and Ian H. Witten. Modeling for Optimal Probability Prediction. ICML. 2002. [[View Context](#)].

C. Titus Brown and Harry W. Bullen and Sean P. Kelly and Robert K. Xiao and Steven G. Satterfield and John G. Hagedorn and Judith E. Devaney. Visualization and Data Mining in an 3D Immersive Environment: Summer Project 2003. [[View Context](#)].

Christos Dimitrakakis and Samy Bengioy. Online Policy Adaptation for Ensemble Classifiers. IDIAP. [[View Context](#)].

Citation Request:

Please refer to the Machine Learning Repository's citation policy

[1] Papers were automatically harvested and associated with this data set, in collaboration with [Rexa.info](#)

Supported By:



In Collaboration With:

