

# Email Spam dataset

Attribute Information for spambase: The last column of “spambase.csv” 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.