

Assignment 3.2

a)

48 continuous values on the range [0, 100] describing the percentage of the following words in the e-mail:

word_freq_make:	continuous.
word_freq_address:	continuous.
word_freq_all:	continuous.
word_freq_3d:	continuous.
word_freq_our:	continuous.
word_freq_over:	continuous.
word_freq_remove:	continuous.
word_freq_internet:	continuous.
word_freq_order:	continuous.
word_freq_mail:	continuous.
word_freq_receive:	continuous.
word_freq_will:	continuous.
word_freq_people:	continuous.
word_freq_report:	continuous.
word_freq_addresses:	continuous.
word_freq_free:	continuous.
word_freq_business:	continuous.
word_freq_email:	continuous.
word_freq_you:	continuous.
word_freq_credit:	continuous.
word_freq_your:	continuous.
word_freq_font:	continuous.
word_freq_000:	continuous.
word_freq_money:	continuous.
word_freq_hp:	continuous.
word_freq_hpl:	continuous.
word_freq_george:	continuous.
word_freq_650:	continuous.
word_freq_lab:	continuous.
word_freq_labs:	continuous.
word_freq_telnet:	continuous.
word_freq_857:	continuous.
word_freq_data:	continuous.
word_freq_415:	continuous.
word_freq_85:	continuous.
word_freq_technology:	continuous.
word_freq_1999:	continuous.
word_freq_parts:	continuous.
word_freq_pm:	continuous.
word_freq_direct:	continuous.
word_freq_cs:	continuous.
word_freq_meeting:	continuous.
word_freq_original:	continuous.
word_freq_project:	continuous.
word_freq_re:	continuous.
word_freq_edu:	continuous.
word_freq_table:	continuous.
word_freq_conference:	continuous.

6 continuous values on the range [0, 100] describing the percentage of the following characters in the e-mail:

char_freq_;	continuous.
char_freq_(:	continuous.
char_freq_[continuous.
char_freq_!	continuous.
char_freq_\$	continuous.
char_freq_#	continuous.

1 continuous real [1,...] value describing the average length of uninterrupted sequences of capital letters

1 continuous integer [1,...] value describing the length of longest uninterrupted sequence of capital letters

1 continuous integer [1,...] value describing the total number of capital letters in the e-mail

b)

Source: https://www.researchgate.net/figure/Top-5-spam-detection-features-19_tbl2_275647490

Example features which could be used to detect spam e-mail:

- Number of capitalized words
- Sum of all character lengths of words
- Number of words containing letters and numbers
- Hour of day when e-mail was sent
- Number of URLs in e-mail