

## CMPE256 - Assignment

### ML Based Spam filter

3. Apply TF and IDF and develop spam filter mode for documents.

Spam Dictionary:

free, click here, visit, open attachment,  
call this number, money, out, extra,  
offer, available, Pension, opportunity,  
chance, investment, Pension.

Solution:  $t_1 = \text{free}$

~~$t_2 = \text{click here}$~~

$t_3 = \text{visit}$

$t_4 = \text{open attachment}$

$t_5 = \text{call this number}$

$t_6 = \text{money}$

$t_7 = \text{out}$

$t_8 = \text{extra}$

$t_9 = \text{offer}$

$t_{10} = \text{available}$

$t_{11} = \text{pension}$

$t_{12} = \text{opportunity}$

$t_{13} = \text{chance}$

$t_{14} = \text{Investment}$

~~Terms~~

~~free~~

~~click here~~

~~visit~~

~~open attachment~~

~~document / term~~

~~Documents~~

	Terms	$d_1$	$d_2$	$d_3$	$d_4$	$d_5$	$d_6$
$t_1$	free	1	1	1	2	0	0
$t_2$	<del>click here</del>	0	0	0	1	0	0
$t_3$	visit	1	1	0	0	0	0
$t_4$	<del>open attachment</del>	0	0	0	0	0	0
$t_5$	<del>call this number</del>	0	0	0	0	0	0
$t_6$	money	0	0	0	0	1	0
$t_7$	out	0	0	1	0	0	0
$t_8$	extra	0	0	0	0	2	0
$t_9$	offer	1	1	0	0	2	4
$t_{10}$	available	1	<del>1</del>	0	0	0	0
$t_{11}$	pension	0	0	2	0	0	0
$t_{12}$	opportunity	0	0	0	0	1	0
$t_{13}$	chance	0	0	0	1	1	0
$t_{14}$	Investment	0	0	0	0	0	1

## Term frequency table:

	$d_1$	$d_2$	$d_3$	$d_4$	$d_5$	$d_6$
$t_1$	1	1	1	1.11	0	0
$t_2$	0	6	0	1	0	9
$t_3$	1	1	0	0	0	0
$t_4$	0	0	0	0	0	0
$t_5$	0	0	0	0	0	0
$t_6$	0	0	0	0	1	0
$t_7$	0	0	1	0	0	0
$t_8$	0	0	0	0	1.11	0
$t_9$	1	1	0	0	1.11	1.204
$t_{10}$	1	1	0	0	0	0
$t_{11}$	0	0	1.11	0	0	0
$t_{12}$	0	0	0	0	1	0
$t_{13}$	0	0	0	1	1	0
$t_{14}$	0	0	0	0	0	1

IDF calculations below.

$$IDF: \cancel{d_1} = 0.2430$$

$$t_2 = 0.845$$

$$t_3 = 0.5440$$

$$t_4 = 0.00$$

$$t_5 = 0.00$$

$$t_6 = 0.845$$

$$t_7 = 0.845$$

$$t_8 = 0.845$$

$$t_9 = 0.243$$

$$t_{10} = 0.5440$$

$$t_{11} = 0.845$$

$$t_{12} = 0.845$$

$$t_{13} = 0.5440$$

$$t_{14} = 0.845$$

\* Tf-IDF calculation below for spam words in all documents.

	$d_1$	$d_2$	$d_3$	$d_4$	$d_5$	$d_6$
$t_1$	0.2430	0.2430	0.2430	0.9379	0	0
$t_2$	0	0	0	0.845	0	0
$t_3$	0.5440	0.5440	0	0	0	0
$t_4$	0	0	0	0	0	0
$t_5$	0	0	0	0	0	0
$t_6$	0	0	0	0	0.845	0
$t_7$	0	0	0.845	0	0	0
$t_8$	0	0	0	0	0.9379	0
$t_9$	0.2430	0.2430	0	0	0.9379	0.2925
$t_{10}$	0.5440	0.5440	0	0	0	0
$t_{11}$	0	0	0.9379	0	0	0
$t_{12}$	0	0	0	0	0.845	0
$t_{13}$	0	0	0	0.5440	0.5440	0
$t_{14}$	0	0	0	0	0	0.845