

TEXT CLASSIFICATION BY NB (group: 404 name not found)

	parliament	italian	loose	win	soccer	championship	government	Class
d1	0	1	0	1	1	1	0	sport
d2	0	0	1	0	1	1	1	sport
d3	0	0	0	0	0	1	0	sport
d4	1	1	1	0	0	0	1	politcs
d5	1	1	0	1	0	0	1	politcs

1. Feature selection by using entropy – select the first k=5 best features.**PARLIAMENT**

$$E(S, \text{parliament}) = -p(\text{sport} | \text{parliament}) * \log_2(\text{sport} | \text{parliament}) - p(\text{politics} | \text{parliament}) * \log_2(\text{politics} | \text{parliament})$$

$$= (-0 * 0) * (-1 * 0) = 0$$

$$p(\text{sport} | \text{parliament}) = 0$$

$$\log_2 p(\text{sport} | \text{parliament}) = 0$$

$$p(\text{politics} | \text{parliament}) = 1$$

$$\log_2 p(\text{politics} | \text{parliament}) = 0$$

ITALIAN

$$E(S, \text{italian}) = -p(\text{sport} | \text{italian}) * \log_2(\text{sport} | \text{italian}) - p(\text{politics} | \text{italian}) * \log_2(\text{politics} | \text{italian}) =$$

$$= (-0.33 * -1.60) * (-0.66 * -0.60) = 0.528 * 0.396 = 0.209088$$

$$p(\text{sport} | \text{italian}) = 0.33$$

$$\log_2 p(\text{sport} | \text{italian}) = -1.60$$

$$p(\text{politics} | \text{italian}) = 0.66$$

$$\log_2 p(\text{politics} | \text{italian}) = -0.60$$

LOOSE

$$E(S, \text{loose}) = -p(\text{sport} | \text{loose}) * \log_2(\text{sport} | \text{loose}) - p(\text{politics} | \text{loose}) * \log_2(\text{politics} | \text{loose}) =$$

$$= (-0.5 * -1) * (-0.5 * -1) = 0.5 * 0.5 = 0.25$$

$$p(\text{sport} | \text{loose}) = 0.5$$

$$\log_2 p(\text{sport} | \text{loose}) = -1$$

$$p(\text{politics} | \text{loose}) = 0.5$$

$$\log_2 p(\text{politics} | \text{loose}) = -1$$

WIN

$$E(S, \text{win}) = -p(\text{sport} | \text{win}) * \log_2(\text{sport} | \text{win}) - p(\text{politics} | \text{win}) * \log_2(\text{politics} | \text{win}) =$$

$$= (-0.5 * -1) * (-0.5 * -1) = 0.5 * 0.5 = 0.25$$

$$p(\text{sport} | \text{win}) = 0.5$$

$$\log_2 p(\text{sport} | \text{win}) = -1$$

$$p(\text{politics} | \text{win}) = 0.5$$

$$\log_2 p(\text{politics} | \text{win}) = -1$$

SOCCER

$$E(S, \text{soccer}) = -p(\text{sport} | \text{soccer}) * \log_2(\text{sport} | \text{soccer}) - p(\text{politics} | \text{soccer}) * \log_2(\text{politics} | \text{soccer}) =$$

$$= (-1 * 0) * (0 * 0) = 0$$

$$p(\text{sport} | \text{soccer}) = 1$$

$$\log_2 p(\text{sport} | \text{soccer}) = 0$$

$$p(\text{politics} | \text{soccer}) = 0$$

$$\log_2 p(\text{politics} | \text{soccer}) = 0$$

CHAMPIONSHIP

$$E(S, \text{championship}) = -p(\text{sport} | \text{championship}) * \log_2(\text{sport} | \text{championship}) - p(\text{politics} | \text{championship}) * \log_2(\text{politics} | \text{championship}) =$$

$$= (-1 * 0) * (0 * 0) = 0$$

$$p(\text{sport} | \text{championship}) = 1$$

$$\log_2 p(\text{sport} | \text{championship}) = 0$$

$$p(\text{politics} | \text{championship}) = 0$$

$$\log_2 p(\text{politics} \mid \text{championship}) = 0$$

GOVERNMENT

$$E(S, \text{government}) = -p(\text{sport} \mid \text{government}) * \log_2(\text{sport} \mid \text{government}) - p(\text{politics} \mid \text{government}) * \log_2(\text{politics} \mid \text{government}) = (0*0)*(-0.66*-0.60) = 0*0.396 = 0$$

$$p(\text{sport} \mid \text{government}) = 0$$

$$\log_2 p(\text{sport} \mid \text{government}) = 0$$

$$p(\text{politics} \mid \text{government}) = 0.66$$

$$\log_2 p(\text{politics} \mid \text{government}) = -0.60$$

So, the first k=5 best features are:

- i. Parliament
- ii. Soccer
- iii. Championship
- iv. Government
- v. italian

2. Prior probabilities estimates, over the (reduced) training set, by using both fractions and m-estimates.

$$P_{c1} = \frac{W_{c1}}{|V_{oc}|} = \frac{3}{5} = 0.6$$

$$P_{c2} = \frac{W_{c2}}{|V_{oc}|} = \frac{3}{5} = 0.6$$

$$P(w \mid c) = p(\text{parliament} \mid \text{sport}) = \frac{d_{w1,c1} + W_{c1}}{d_{c1} + |V_{oc}|} = \frac{0+3}{3+5} = 0.375$$

$$P(w \mid c) = p(\text{parliament} \mid \text{politics}) = \frac{d_{w1,c2} + W_{c2}}{d_{c2} + |V_{oc}|} = \frac{2+3}{3+5} = 0.625$$

$$P(w \mid c) = p(\text{soccer} \mid \text{sport}) = \frac{d_{w2,c1} + W_{c1}}{d_{c1} + |V_{oc}|} = \frac{2+3}{3+5} = 0.625$$

$$P(w \mid c) = p(\text{soccer} \mid \text{politics}) = \frac{d_{w2,c2} + W_{c2}}{d_{c2} + |V_{oc}|} = \frac{0+3}{3+5} = 0.375$$

$$P(w \mid c) = p(\text{championship} \mid \text{sport}) = \frac{d_{w3,c1} + W_{c1}}{d_{c1} + |V_{oc}|} = \frac{3+3}{3+5} = 0.75$$

$$P(w \mid c) = p(\text{championship} \mid \text{politics}) = \frac{d_{w3,c2} + W_{c2}}{d_{c2} + |V_{oc}|} = \frac{0+3}{3+5} = 0.375$$

$$P(w \mid c) = p(\text{government} \mid \text{sport}) = \frac{d_{w4,c1} + W_{c1}}{d_{c1} + |V_{oc}|} = \frac{1+3}{3+5} = 0.5$$

$$P(w \mid c) = p(\text{government} \mid \text{politics}) = \frac{d_{w4,c2} + W_{c2}}{d_{c2} + |V_{oc}|} = \frac{2+3}{3+5} = 0.625$$

$$P(w \mid c) = p(\text{italian} \mid \text{sport}) = \frac{d_{w5,c1} + W_{c1}}{d_{c1} + |V_{oc}|} = \frac{1+3}{3+5} = 0.5$$

$$P(w \mid c) = p(\text{italian} \mid \text{politics}) = \frac{d_{w5,c2} + W_{c2}}{d_{c2} + |V_{oc}|} = \frac{2+3}{3+5} = 0.625$$

3. Classification of the following docs:

- o d1 = "The team of Juventus will likely not win the Italian soccer championship".

	Parliament	Italian	Soccer	Championship	Government	CLASS
d1	0	1	1	1	0	SPORT

$$p(\text{sport} \mid \{\text{italian}, \text{soccer}, \text{championship}\}) = p(\text{sport}) * p(\text{italian} \mid \text{sport}) * p(\text{soccer} \mid \text{sport}) * p(\text{championship} \mid \text{sport}) = 0.6 * 0.5 * 0.625 * 0.75 = \mathbf{0.14}$$

$$p(\text{politics} \mid \{\text{italian}, \text{soccer}, \text{championship}\}) = p(\text{politics}) * p(\text{italian} \mid \text{politics}) * p(\text{soccer} \mid \text{politics}) * p(\text{championship} \mid \text{politics}) = 0.4 * 0.625 * 0.375 * 0.375 = \mathbf{0.04}$$

- o d2 = "the Italian parliament, based on the initiative of the government, has approved a law for Italian soccer teams".

	Parliament	Italian	Soccer	Championship	Government	CLASS
d2	1	1	1	0	1	POLITICS

$$p(\text{sport} \mid \{\text{parliament}, \text{italian}, \text{soccer}, \text{government}\}) = p(\text{sport}) * p(\text{parliament} \mid \text{sport}) * p(\text{italian} \mid \text{sport}) * p(\text{soccer} \mid \text{sport}) * p(\text{government} \mid \text{sport}) = 0.6 * 0.375 * 0.5 * 0.625 * 0.5 = \mathbf{0.036}$$

$$p(\text{politics} \mid \{\text{parliament}, \text{italian}, \text{soccer}, \text{government}\}) = p(\text{politics}) * p(\text{parliament} \mid \text{politics}) * p(\text{italian} \mid \text{politics}) * p(\text{soccer} \mid \text{politics}) * p(\text{government} \mid \text{politics}) = 0.4 * 0.625 * 0.625 * 0.375 * 0.625 = \mathbf{0.037}$$