

1. Naive Bayes text classification

(i) Bernoulli model

Document Term Matrix: (Binary)

	Hualien	Taiwan	Sapporo	Shanghai	Japan	class
d ₁	1	1	0	0	0	Y
d ₂	1	1	0	0	0	Y
d ₃	0	0	1	1	0	Y
d ₄	1	0	1	0	1	N
d ₅	0	1	1	0	1	N
d ₆	1	1	0	0	0	?

$$\text{class: YES} \Rightarrow \frac{3}{5} \times \frac{2+1}{3+2} \times \frac{2+1}{3+2} \times \frac{2+1}{3+2} \times \frac{2+1}{3+2} \times \frac{3+1}{3+2} = 0.062 \dots$$

$$\text{class: No} \Rightarrow \frac{2}{5} \times \frac{1+1}{2+2} \times \frac{1+1}{2+2} \times \frac{0+1}{2+2} \times \frac{2+1}{2+2} \times \frac{0+1}{2+2} = 0.004 \dots$$

Ans: class label: YES

(ii) Multinomial model

Label Term Matrix: (count)

	Hualien	Taiwan	Sapporo	Shanghai	Japan
YES	3	3	1	1	0
NO	1	1	2	0	3

$$\text{class: YES} \Rightarrow \frac{3}{5} \times \frac{3+1}{8+5} \times \frac{3+1}{8+5} \times \frac{3+1}{8+5} = 0.017 \dots$$

$$\text{class: No} \Rightarrow \frac{2}{5} \times \frac{1+1}{7+5} \times \frac{1+1}{8+5} \times \frac{1+1}{8+5} = 0.001 \dots$$

Ans: class label: YES

2.

(i)

$$M_1 = \begin{array}{c|cccc} & \textcircled{1} & \textcircled{2} & \textcircled{3} & \textcircled{4} \\ \hline \textcircled{1} & 0 & \frac{1}{2} & \frac{1}{2} & 0 \\ \textcircled{2} & \frac{1}{2} & 0 & 0 & \frac{1}{2} \\ \textcircled{3} & \frac{1}{4} & \frac{1}{4} & \frac{1}{4} & \frac{1}{4} \\ \textcircled{4} & \frac{1}{3} & \frac{1}{3} & \frac{1}{3} & 0 \end{array}$$

 $M_2 :$

$$\begin{array}{c|cccc} & \textcircled{1} & \textcircled{2} & \textcircled{3} & \textcircled{4} \\ \hline \textcircled{1} & \frac{1}{4} & \frac{1}{4} & \frac{1}{4} & \frac{1}{4} \\ \textcircled{2} & \frac{1}{4} & \frac{1}{4} & \frac{1}{4} & \frac{1}{4} \\ \textcircled{3} & \frac{1}{4} & \frac{1}{4} & \frac{1}{4} & \frac{1}{4} \\ \textcircled{4} & \frac{1}{4} & \frac{1}{4} & \frac{1}{4} & \frac{1}{4} \end{array}$$

$$\underline{0.8 M_1 + 0.2 M_2 =}$$

$$\begin{array}{c|cccc} & \textcircled{1} & \textcircled{2} & \textcircled{3} & \textcircled{4} \\ \hline \textcircled{1} & \frac{1}{20} & \frac{9}{20} & \frac{9}{20} & \frac{1}{20} \\ \textcircled{2} & \frac{9}{20} & \frac{1}{20} & \frac{1}{20} & \frac{9}{20} \\ \textcircled{3} & \frac{1}{4} & \frac{1}{4} & \frac{1}{4} & \frac{1}{4} \\ \textcircled{4} & \frac{19}{60} & \frac{19}{60} & \frac{19}{60} & \frac{1}{20} \end{array}$$

#

(ii)

$$\text{step 0 : } \langle 1, 0, 0, 0 \rangle$$

$$\text{step 1 : } \langle \frac{1}{20}, \frac{9}{20}, \frac{9}{20}, \frac{1}{20} \rangle$$

$$\begin{aligned} \text{step 2 : } & \langle \frac{1}{20} \times \frac{1}{20} + \frac{9}{20} \times \frac{9}{20} + \frac{9}{20} \times \frac{1}{4} + \frac{1}{20} \times \frac{19}{60}, \\ & \frac{1}{20} \times \frac{9}{20} + \frac{9}{20} \times \frac{1}{20} + \frac{9}{20} \times \frac{1}{4} + \frac{1}{20} \times \frac{19}{60}, \\ & \frac{1}{20} \times \frac{9}{20} + \frac{9}{20} \times \frac{1}{20} + \frac{9}{20} \times \frac{1}{4} + \frac{1}{20} \times \frac{19}{60}, \\ & \frac{1}{20} \times \frac{1}{20} + \frac{9}{20} \times \frac{9}{20} + \frac{9}{20} \times \frac{1}{4} + \frac{1}{20} \times \frac{1}{20} \rangle \end{aligned}$$

(iii)

$$\langle 1, 0, 0, 0 \rangle \times M^{10}$$

3.

$$n = 20$$

$$\text{Total number of pairs} = 20 \times 19 / 2 = 190$$

$$TP + FP = C_2^7 + C_2^8 + C_2^5 = 59$$

$$TP = C_2^4 + C_2^3 + C_2^5 + C_2^2 + C_2^3 + C_2^2 = 24$$

$$FP = 59 - 24 = 35$$

$$FN + TN = 190 - 59 = 131$$

$$FN = (4 \times 3 + 1 \times 2) + (3 \times 5) + (2 \times 3) = 35$$

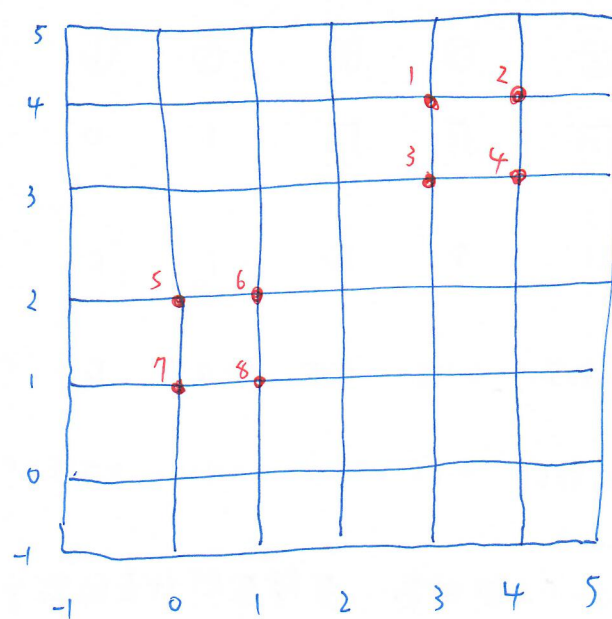
$$TN = 131 - 35 = 96$$

Confusion Matrix :

24		35	
	TP	FN	
35		96	
	FP	TN	

$$\text{Rand Index} : \frac{24 + 96}{24 + 35 + 35 + 96} = \frac{120}{190} = 0.63 \quad \#$$

4. (i) My idea: 使用 PCA 降維, 並人工選取起始點。



example: 當分 2 群時, 可以選 2, 7 當作 seed。

(ii) K-means ++

step 1: Randomly select the first centroid from the data points.

step 2: For each data point compute its distance from the nearest, previously chosen centroid.

step 3: Select the next centroid from the data points such that the probability of choosing a point as centroid is directly proportional to its distance from the nearest, previously chosen centroid.

step 4: Repeat steps 2 and step 3 until K centroids have been sampled.

example: 以 step 1 選 6 號為例:

	①	②	③	④	⑤	⑥	⑦	⑧
$D(x)$	$2\sqrt{2}$	$\sqrt{13}$	$\sqrt{5}$	$\sqrt{10}$	1	0	$\sqrt{2}$	1
$D(x)^2$	8	13	5	10	1	0	2	1
$P(x)$	0.2	0.325	0.125	0.25	0.025	0	0.05	0.025
sum	0.2	0.525	0.65	0.9	0.925	0.925	0.975	1

隨機產生一個 0-1 之間的數, 選擇區間代表的點當作 seed.

ex: ①的區間為 0-0.2, ②的區間為 0.2-0.525。

而 ①②③④ 就佔了整體的 90%, 確實與現有的 ⑥ 較遠。