

SEE Regression problems

classmate
Date _____
Page _____

*

$$y = mx + c$$

$$c = \frac{\sum y (\sum x^2) - \sum x \sum xy}{n(\sum x^2) - (\sum x)^2}$$

$$m = \frac{n \sum xy - \sum x \sum y}{n(\sum x^2) - (\sum x)^2}$$

x	y	xy	x^2	Predicted	Error
1	2	2	1	2.2	0.8
2	4	8	4	3.4	0.6
3	5	15	9	4	1
4	4	16	16	4.6	0.6
5	5	25	25	5.2	0.2
Σ	15	20	66	55	

$$c = \frac{20 \times 255 - 15 \times 66}{5 \times 255 - 225}$$

$$c = \underline{\underline{2.2}}$$

$$m = \frac{5 \times 66 - 15 \times 20}{5 \times 55 - 225}$$

$$m = \underline{\underline{0.6}}$$

$$y = 0.6x + 2.2$$

Q If the weather is sunny, the player should play or not

- Frequency Table for Weather Conditions

Outlook	Yes	No
Rainy	2	2
Overcast	5	0
Sunny	3	2
Total	10	4

frequency table for all independent variables (here only)

Likelihood Table

outlook	Probability	Yes	
		No	Yes
Rainy	$4/14 = 0.29$	2	2
overcast	$5/14 = 0.35$	5	0
Sunny	$5/14 = 0.35$	3	2
All		$10/14 = 0.71$	$4/14 = 0.29$

Applying Bayes Theorem

$$P(\text{Yes}|\text{Sunny}) = \frac{P(\text{Sunny}|\text{Yes}) \times P(\text{Yes})}{P(\text{Sunny})}$$

$$P(\text{Sunny}|\text{Yes}) = \frac{3}{10} = 0.3$$

$$P(\text{Sunny}) = 0.35$$

$$P(\text{Yes}) = 0.71$$

$$P(\text{Yes}|\text{Sunny}) = \frac{0.3 \times 0.71}{0.35} = 0.608$$

$$P(\text{No/Sunny}) = \frac{P(\text{Sunny/No}) * P(\text{No})}{P(\text{Sunny})}$$

$$P(\text{Sunny/No}) = \frac{2}{4} = \underline{\underline{0.5}}$$

$$P(\text{No}) = 0.29$$

$$P(\text{Sunny}) = 0.35$$

$$P(\text{No/Sunny}) = \frac{0.5 * 0.29}{0.35}$$

$$= \underline{\underline{0.41}}$$

From the above calculation

$$P(\text{Yes/Sunny}) > P(\text{No/Sunny})$$

Hence on a sunny day, player can play the game

Q.	No.	Colour	Legs	Height	Smelly	Species
	1	White	3	Short	Yes	M
	2	Green	2	Tall	No	M
	3	Green	3	Short	Yes	M
	4	White	3	Short	Yes	M
	5	Green	2	Short	No	H
	6	White	2	Tall	No	H
	7	White	2	Tall	No	H
	8	White	2	Short	Yes	H

New instance: If colour = green, legs = tall & smelly = No, identify which species it belongs to

$$P(M) = \frac{4}{8} = 0.5 ; P(H) = 0.5$$

Likelihood Frequency table for colour

Colour	M(4)	H(4)	Probability
White	2/4	3/4	5/8 = 0.625
Green	2/4	1/4	3/8 = 0.375
	4/8 = 0.5	4/8 = 0.5	

→ for Legs

Legs	M(u)	H(u)	Prob
3	3/4	0	3/8 = 0.375
2	1/4	4/4	5/8 = 0.625
	4/8 = 0.5	4/8 = 0.5	

→ Height

	M	H	Prob
Short	3/4	2/4	5/8
Tall	1/4	2/4	3/8

→ Smelly

	M	H	Prob
Yes	3/4	1/4	
No	1/4	3/4	

0.0117

0.0156

CLASSEmate

Date _____
Page _____

New instance = (Colour = green, Legs = 2, Height = Tall, Smelly = No)

$P(M | \text{New_instance})$

$$= P(M) * P(\text{Color} = \text{green} | M) * P(\text{Legs} = 2 | M) * P(\text{Height} = \text{Tall} | M) * P(\text{Smelly} = \text{No} | M)$$

$P(H | \text{New_instance})$

$$= 0.5 * \frac{2}{4} * \frac{1}{4} * \frac{1}{4} * \frac{1}{4} = \underline{\underline{0.0156}} = 0.0039$$

$P(H | \text{New_instance})$

$$= P(H) * P(\text{green} | H) * P(2 | H) * P(\text{Tall} | H) * P(\text{No} | H)$$

$$= 0.5 * 0.25 * 1 * 0.5 * 0.75$$

$$= \underline{\underline{0.0468}}$$

$P(H) > P(M)$

$P(H | \text{New_instance}) > P(M | \text{New_instance})$

$\therefore \text{Species} = H$

* Newinst = (Colour = green, Legs = 2, Height = Tall, Smelly = Yes)

$P(M | \text{Newinst})$

$$= P(M) * P(\text{green} | M) * P(2 | M) * P(\text{Tall} | M) * P(\text{Yes} | M)$$

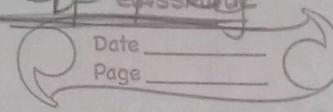
$$= 0.5 *$$

$\underline{\underline{\text{P}(H | \text{Species} = H)}}$

Even though Species is M in dataset, there is chance to get H in naive bayes

Q. Cluster Analysis with the help of ~~classmate~~

V.I.M.P



Date _____

Page _____

- * we can't classify based on shape, size, color. e.g. apple, orange & ball are all round but don't fall under same cluster

Cluster is formed based on similarity, distance & centroid

→ Problems with Measuring Distances

→ A more general formula for distances

$$d_m(x, y) = \sqrt{\sum_{i=1}^n d(x_i, y_i)}$$

* for continuous attributes $d(x_i, y_i) = (x_i - y_i)^2$

* for discrete

V.I.M.P Which cluster should an example belong to?

* $d(x, c_i)$ → distance of x from centroid

e.g.: $c_1 = (3, 4)$, $c_2 = (4, 6)$, $c_3 = (5, 7)$

$$x = (4, 4)$$

$$d(x, c_1) = \sqrt{1^2 + 0^2} = 1$$

$$d(x, c_2) = \sqrt{0^2 + 2^2} = 2$$

$$d(x, c_3) = \sqrt{10} = 3$$

here the shortest distance is $d(x, c_1)$

∴ x belongs to cluster c_1

* If there are 10 datapoints initially! then there are 10 clusters, this can be reduced by using K-means clustering!

Q. Hours Study	Pass(1) / Fail(0)
29	0
15	0
33	1
28	0
39	1

Use Logistic regression as classifier to answer the following questions

1. Calculate the probability for the student who studied 33 hrs. given $\log(\text{odds}) = -64 + 2 \times \text{hours}$

$$P = \frac{1}{1 + e^{-x}} \quad S(x) = \frac{1}{1 + e^{-x}}$$

$$x = -64 + 2 \times 33 = -64 + 66 = 2$$

$$P = \frac{1}{1 + e^{-x}} = \frac{0.88}{-}$$

There is 88% chance that he will pass

2. At least how many hours student should study that makes he will pass the course with prob of more than 95%.

$$P = \frac{1}{1 + e^{-x}} = 0.95 \quad \underline{\underline{0.95}}$$

$$0.95 * (1 + e^{-x}) = 1$$

$$0.95 * e^{-x} = 1 - 0.95$$

classmate
Date _____
Page _____

19-3-24

$$Z = 2.95$$

$$\text{log(odds)} = Z = -64 + 2 * \text{hours}$$

$$2.94 = -64 + 2 * \text{hours}$$

$$2 * \text{hours} = 2.94 + 64 = 66.94$$

$$\text{hours} = \frac{66.94}{2} = \underline{\underline{33.47 \text{ hours}}}$$

K-M
* firm
such
cluster
* item
(other)

Example: The table below

	Group-1	Group-2	Group-3
(2, 5)	(4, 3)	(1, 5)	
(1, 4)	(3, 1)	(3, 1)	
(3, 6)	(2, 2)	(2, 3)	
centroids	(2, 5)	(3, 4)	(2, 3)

* $(2, 5) (2, 5)$ $d = 0 \rightarrow$ least dist. \Rightarrow both clusters

$$(2, 5) (3, 4) \quad d = \sqrt{1+1} = \underline{\underline{\sqrt{2}}}$$

$$(2, 5) (2, 3) \quad d = \cancel{\sqrt{2}} = \underline{\underline{2}}$$

* $(1, 4) (2, 5)$ $d = \sqrt{2} \Rightarrow \min \Rightarrow$ right cluster

$$(1, 4) (3, 4) \quad d = \sqrt{2+0} = \sqrt{2^2} = 2$$

$$(1, 4) (2, 3) \quad d = \sqrt{2}$$

* $(3, 6) (2, 5)$ $d = \sqrt{2} \Rightarrow \min \Rightarrow$ right cluster

$$(3, 6) (3, 4) \quad d = \sqrt{2^2} = 2$$

$$(3, 6) (2, 3) \quad d = \sqrt{1^2 + 3^2} = \sqrt{10}$$

Group-2

* $(4, 3) (2, 5)$ $d = \sqrt{2^2 + 2^2} = \sqrt{8}$

$$(4, 3) (3, 4) \quad d = \sqrt{1^2 + 1^2} = \sqrt{2}$$

$$(4, 3) (2, 3) \quad d = \sqrt{2^2} = 2$$

#~~Aug, 23rd~~

* $(3, 1)(2, 5)$

$$d = \sqrt{1^2 + 2^2} = \sqrt{5} \rightarrow \text{min}$$

$(3, 1)(3, 4)$

$$d = \sqrt{0 + 3^2} = 3 \quad \text{wrong}$$

$(3, 1)(2, 3)$

$$d = \sqrt{1^2 + 4^2} = \sqrt{17} \quad \text{cluster}$$

* $(2, 2)(2, 5)$

$$d = \sqrt{0 + 3^2} = 3$$

$(2, 2)(3, 4)$

$$d = \sqrt{1 + 2^2} = \sqrt{5}$$

$(2, 2)(2, 3)$

$$d = \sqrt{0 + 1^2} = 1 \rightarrow \text{min}$$

Group 3

* $(1, 5)(2, 5)$ $d = \sqrt{1^2 + 0} = 1 \rightarrow \text{min}$ ✗

$(1, 5)(3, 4)$ $d = \sqrt{2^2 + 1^2} = \sqrt{5}$

$(1, 5)(2, 3)$ $d = \sqrt{1 + 2^2} = \sqrt{5}$

* $(3, 1)(2, 5)$ $d = \sqrt{1^2 + 4^2} = \sqrt{17}$

$(3, 1)(3, 4)$ $d = \sqrt{0 + 3^2} = 3$

$(3, 1)(2, 3)$ $d = \sqrt{1 + 2^2} = \sqrt{5} \rightarrow \text{min}$ ✓

* $(2, 3)(2, 5)$

$(2, 3)(3, 4)$

$(2, 3)(2, 3)$

Group 1	Group 2	Group 3
(2, 5)	(4, 3)	(2, 2)
(1, 4)		(3, 1)
(3, 6)		(2, 3)
(3, 7)		
(1, 5)		
centroids	(2, 5.4)	(2.3, 2)
	(4, 3)	

Q. Apply K-means algorithm to solve the given eg into 3 clusters

	x_i	y_i	(2, 10)	(5, 8)	(1, 2)	cluster	(2, 10)	(6, 1)
A1	2	10	0	$\sqrt{3^2+2^2}$ 3.61	8.06	1	0	$\sqrt{5^2+6^2}$ 5.6
A2	2	5	5	4.24	3.16	3	5	$\sqrt{4^2+1^2}$ $\sqrt{17}$
A3	8	4	8.49	5	7.28	2		
B1	5	8	$\sqrt{3^2+5^2}$ 3.61	0	7.21	2		
B2	1	5	$\sqrt{5^2+5^2}$ 7.07	3.61	6.11	2		
B3	6	4	$\sqrt{4^2+6^2}$ 7.21	4.12	5.39	2		
C1	1	2	$\sqrt{1^2+8^2}$ 8.06	7.21	0	3		
C2	4	9	$\sqrt{2^2+1^2}$ 2.24	1.41	7.62	2		

centroids: (2, 10) (5, 8) (1, 2)

Update centroids

$$\frac{8+5+7+6+4}{5} = 6$$

$$\frac{4+8+5+4+9}{5} = 6$$

$$\frac{2+1}{2} = 3$$

(New)

Updated centroids

cluster	(2, 10)	(6, 6)	(1.5, 3.5)	New Cluster	(3, 9.5)	(6.5, 5.25)	(1.5, 3.5)
0	$\sqrt{4^2 + 6^2}$ 5	$\sqrt{5^2 + 6^2}$ $\sqrt{65}$	0.77		1		
5				3			
2				2			
2				2	2.5	3.13	4.87
2				2			
3			2.5	3			
3			3	1			
2				2			

New Clusters: 1, 3, 2, 4, 2, 2, 3, 4
new

$$\frac{2+4}{2} = 3; \quad \frac{10+9}{2} = 9.5 \quad | \quad 8+5+7+6$$

$x_1 \mid y_1$

2 | 10

2 | 5

$A_1, B_1, C_2 \Rightarrow \text{Cluster 1}$

$A_3, B_2, B_3 \Rightarrow \text{Cluster 2}$

$A_2, C_1 \Rightarrow \text{Cluster 3}$

20-3-24

Q.	X	Y	(182, 72)	(180, 56)	
③ 168	60		18.43	4.48	$\rightarrow \text{Cluster 2}$
④ 179	68				$\rightarrow 1$
⑤ 182	72				$\rightarrow 1$
⑥ 188	71				$\rightarrow 1$
⑦ 180	71				$\rightarrow 1$
⑧ 180	70				$\rightarrow 1$
⑨ 183	84				$\rightarrow 1$
⑩ 180	88				$\rightarrow 1$
⑪ 180	61				$\rightarrow 1$
⑫ 171	76		21.18	21.18	$\rightarrow 1$
Centroid		185	72		$\rightarrow 1$
Centroid		170	56		$\rightarrow 2$

cluster
 (179, 61)
 (182, 72)
 (188, 56)
 (180, 56)
 (180, 43)
 (183, 56)
 (18, 56)
 (171, 56)
 (159, 56)

Updated
Centroids: (1

Cluster 1	Cluster 2
(179, 68)	(170, 56)
(182, 72)	(168, 60)
(188, 71)	
(180, 71)	
(180, 70)	
(183, 84)	
(180, 88)	
(180, 67)	
(177, 76)	
(185, 72)	
Updated Centroids: (181.4, 74.5)	(169, 58)

K1: 1, 4, 5, 6, 7, 8, 9, 10, 11, 12

K2: 2, 3

A

$$x_1 = (1, 0)$$

$$x_2 = (2, 2)$$

B

$$y_1 = (3, 3)$$

$$y_2 = (4, 4)$$

New points: $(1, 3), (0, 3)$

$$d(x_1, y_1) = \sqrt{(3-1)^2 + (3-0)^2} = \sqrt{4+9} = \sqrt{13}$$

$$d(x_1, y_2) = \sqrt{(4-1)^2 + (4-0)^2} = \sqrt{3^2 + 4^2} = \sqrt{25} = 5$$

$$d(x_2, y_1) = \sqrt{(3-2)^2 + (3-2)^2} = \sqrt{1+1} = \sqrt{2}$$

$$d(x_2, y_2) = \sqrt{(4-2)^2 + (4-2)^2} = \sqrt{4+4} = \sqrt{8}$$

The smallest of these values is $d(x_2, y_1) = \sqrt{2}$
∴ distance b/w 2 clusters is $\underline{\underline{d(A, B) = \sqrt{2}}}$

IMP (3m)

Normalization

$$x_i = \frac{x_i}{\sum x_j}$$

Q. Consider the following set of 6 one-dimensional ~~set~~ of data points:

18, 22, 25, 42, 27, 43.

Apply agglomerative hierarchical algo to build a dendrogram. Merge the clusters using minimum distance & update the proximity matrix accordingly

	18	22	25	27	42	43
18	0	4	7	9	24	25
22	4	0	3	5	20	21
25	7	3	0	2	17	18
27	9	5	2	0	15	16
42	24	20	17	15	0	1
43	25	21	18	16	1	0

* Identify the min (cost) value other than 0 $\Rightarrow (42, 43)$

Grouping rows

	18	22	25	27	(42, 43)
18	0	4	7	9	24
22	4	0	3	5	20
25	7	3	0	2	17
27	9	5	2	0	15
(42, 43)	24	20	17	15	(25, 27)

(42, 43) (25, 27)

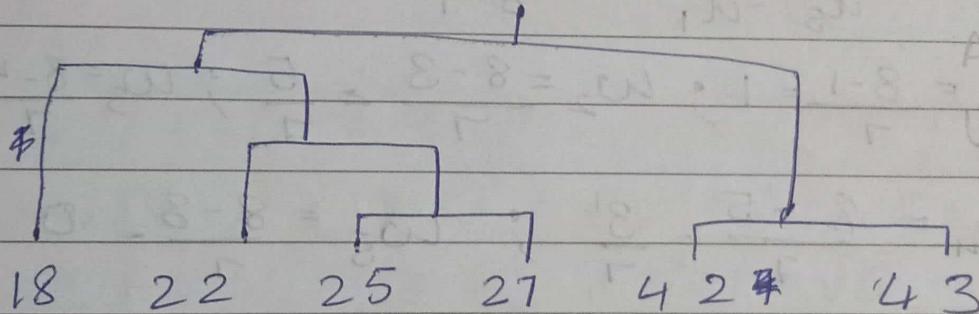
	18	22	(25, 27)	(42, 43)
18	0	4	7	24
22	4	0	3	20
(25, 27)	7	3	0	15
(42, 43)	24	20	15	0

(42, 43) ((25, 27), 22)

	18	(22, 25, 27)	(42, 43)
18	0	4	24
(22, 25, 27)	4	0	15
(42, 43)	24	15	0

((42, 43), ((25, 27), 22), 18))

	(18, 22, 25, 27)	(18, 22, 25, 27)	(42, 43)
(18, 22, 25, 27)	0	15	
(42, 43)	15	0	



2-11-24

Weighted Nearest Neighbors

- * we are going to assign weights to each of the neighbors
- * any neighbor nearest to new instance \rightarrow highest weight
- * K neighbors ordered according to their distances d_1, d_2, \dots, d_K from x^* so that d_1 is the smallest & d_K is the greatest

$$w_i = \frac{d_K - d_i}{d_K - d_1}$$

15-4-24

classmate

Date _____
Page _____

Naive Bayes

Q. The test returns a correct +ve result in only 98% of the cases in which the disease is actually present, and a correct -ve result in only 97% of the cases in which the disease is not present. Furthermore, 0.001 of the entire population have this cancer. Does this patient have cancer or not? How does probability of cancer being +ve compared to probability of not having cancer being +ve.

$$P(\text{cancer}) = 0.001$$

$$P(\neg \text{cancer}) = 1 - P(\text{cancer}) = 1 - 0.001 = \underline{\underline{0.999}}$$

$$P(+|\text{cancer}) = 0.98, P(-|\text{cancer}) = 0.02$$

$$P(-|\neg \text{cancer}) = 0.97, P(+|\neg \text{cancer}) = 0.03$$

$P(\text{cancer} | +)$ compare with $P(\neg \text{cancer} | +)$

$$P(A|B) = \frac{P(B|A) \cdot P(A)}{P(B)}$$

$$P(\text{cancer} | +) = \frac{P(+|\text{cancer}) \cdot P(\text{cancer})}{P(+)}$$

$$\begin{aligned}
 &= \frac{0.98 \times 0.001}{P(+)} \rightarrow \text{all +ve cases} \\
 &= \frac{0.98 \times 0.001}{P(+/\text{cancer}) + P(+/\neg\text{cancer})} \\
 &= \frac{0.98 \times 0.001}{P(+/\text{cancer}) P(\text{cancer}) + P(+/\neg\text{cancer}) P(\neg\text{cancer})} \\
 &= \frac{0.98 \times 0.001}{(0.98 \times 0.001) + (0.03 \times 0.999)} \\
 &= \underline{0.031664} = \underline{3.16\%}
 \end{aligned}$$

$P(\text{cancer}/+)$

$$= \frac{P(+/\neg\text{cancer}) \cdot P(\text{cancer})}{P(+)}$$

$$= \frac{0.03 \times 0.001}{P(+)} \times 0.999$$

$$= \frac{0.03 \times 0.999}{P(+, \text{cancer}) + P(+, \neg\text{cancer})}$$

$$= \underline{0.02997}$$

$$= \underline{0.9683}$$

18-4-24

classmate

Date _____
Page _____Agglomerative 2-D Problem

Sample	X	Y	
P ₁	0.40	0.53	0.53
P ₂	0.22	0.38	
P ₃	0.35	0.32	
P ₄	0.26	0.19	
P ₅	0.08	0.41	
P ₆	0.45	0.30	

	P ₁	P ₂	P ₃	P ₄	P ₅	P ₆
P ₁	0					
P ₂	0.23	0				
P ₃	0.22	0.14	0			
P ₄	0.37	0.19	0.13	0		
P ₅	0.34	0.14	0.28	0.23	0	
P ₆	0.24	0.24	0.10	0.22	0.39	0

↓
lowest

(P₃, P₆)

	P ₁	P ₂	P ₃ , P ₆	P ₄	P ₅
P ₁	0				
P ₂	0.23	0			
P ₃ , P ₆	0.22	0.14	0		
P ₄	0.37	0.19	0.13	0	
P ₅	0.34	0.14	0.28	0.23	0

$((P_3, P_6), P_4)$

	P_1	P_2	P_3, P_4, P_6	P_5
P_1	0			
P_2	0.23	0		
P_3, P_6	0.22	0.14	0	
P_5	0.34	0.14	0.23	0

 $((P_3, \cancel{P_6}) P_4), (P_2, P_5)$

	P_1	P_2, P_5	P_3, P_4, P_6
P_1	0		
P_2, P_5	0.23	0	
P_3, P_4, P_6	0.22	0.14	0

 $[((((P_3, P_6), P_4), (P_2, 5)), P_1)]$
