

CS 541: Machine Learning, Spring 2024

Name: SRIHARI TANMAY KARTHIK TADALA

PSU ID: 918597835

Graduate Student

Q

MACHINE LEARNING ASSIGNMENT-3

- A1.) Given us Monty Hall case of 4 doors (3 goats, 1 car, host opens 1 door after selection)

Host knows what is behind each door

Assumption: uniform prior distribution, meaning each door has an equal chance of car initially ie $(1/4)^{\text{th}}$ for each door

Prior probability: $- P(h)$

$$P(h_1) = 1/4, P(h_2) = 1/4, P(h_3) = 1/4, P(h_4) = 1/4$$

Calculate the conditional probabilities:

Door 1 chosen; Monty opens door 2

$$P(D|h_1) = 1/3 \text{ (since Monty has 3 choices)}$$

Door 2 chosen; Monty opens door 1

$$P(D|h_2) = 0 \text{ (Host can't open the chosen door).}$$

Door 3 chosen; Monty opens door 4.

$$P(D|h_3) = 1/2 \text{ (Monty opens and has 2 doors left)}$$

Door 4 chosen; Monty opens door 3

$$P(D|h_4) = 1/2 \text{ (Monty has 2 choices)}$$

Calculating marginal probability $P(D)$:

$$\begin{aligned} P(D) &: [P(D/h_1) \cdot P(h_1) + P(D/h_2) \cdot P(h_2) + P(D/h_3) \cdot P(h_3) \\ &\quad + P(D/h_4) \cdot P(h_4)] \\ &= 1/3(1/4) + 0 + (1/2)(1/4) + (1/2)(1/4) \\ &= 1/12 + 1/8 = 1/3 \end{aligned}$$

$$P(D) = 1/3$$

Calculating posterior probability $P(h_i|D)$:

$$P(h_i|D) = \frac{P(D|h_i) \cdot P(h_i)}{P(D)}$$

$$P(h_1|D) = \frac{1/3 \cdot 1/4}{1/3} = 1/4 \quad P(h_2|D) = 0$$

$$P(h_3|D) = \frac{1/2 \cdot 1/4}{1/3} = 3/8$$

$$P(h_4|D) = \frac{1/2 \cdot 1/4}{1/3} = 3/8$$

The probability of car behind the 2 doors if you switch is $3/8$ which is greater than the probability of first door ($1/4$)

⑦

Therefore it is advantageous to switch initial selection.

2) Example of bias: using the logistic regression when the data is not linearly separable.

Example of variance: - Using SVM; on high dimension feature space can exhibit high variance.

Example of Noise: → Models like Naive Bayes classifier are more sensitive to noisy data

because naive baye classifiers assume independence between features & may struggle with noisy or correlated data.

3.a) To show how new patient P7 is classified under which class)?

Using naive Bayes

1) Use training data to compute probabilistic model.

$$P(\text{cough} = \text{True} | \text{Flu}) = 2/4 \quad \left| \begin{array}{l} P(\text{cough} = \text{True} | \text{cold}) = 1/2 \\ P(\text{cough} = \text{False} | \text{cold}) = 1/2 \end{array} \right.$$

$$P(\text{cough} = \text{False} | \text{Flu}) = 2/4$$

$$P(\text{Fever} = \text{True} | \text{Flu}) = 3/4 \quad \left| \begin{array}{l} P(\text{Fever} = \text{True} | \text{cold}) = 1/2 \\ P(\text{Fever} = \text{False} | \text{cold}) = 1/2 \end{array} \right.$$

$$\begin{array}{ll}
 P(\text{Fever} = \text{False} | \text{Flu}) = 1/4 & P(\text{Fever} = \text{False} | \text{cold}) = 1/2 \\
 P(\text{chills} = \text{True} | \text{Flu}) = 3/4 & P(\text{chills} = \text{True} | \text{cold}) = 1/2 \\
 P(\text{chills} = \text{False} | \text{Flu}) = 1/4 & P(\text{chills} = \text{False} | \text{cold}) = 1/2 \\
 \\
 P(\text{congestion} = \text{True} | \text{False}) = 1/4 & P(\text{congestion} = \text{True} | \text{cold}) = 3/2 \\
 P(\text{congestion} = \text{False} | \text{flu}) = 3/4 & P(\text{congestion} = \text{False} | \text{cold}) = 0
 \end{array}$$

② Let's find the prior class of

$$P(\text{Flu}) = 4/6, P(\text{cold}) = 2/6.$$

③ The probability of the class & cold class for patient P_7 is

$$\begin{aligned}
 \text{Flu class}(x) &= \arg \max_i P(\text{Flu}) \leq P(x_i | \text{Flu}) \cdot \\
 &\quad \times P(\text{Fever} = \text{False} | \text{Flu}) \times P(\text{chills} = \text{False} | \text{Flu}) \times \\
 &\quad P(\text{congestion} = \text{True} | \text{False})
 \end{aligned}$$

$$\text{Flu}(x) = 4/6 \times 2/4 \times 1/4 \times 1/4 \times 1/4$$

$$\text{Flu}(x) = 0.0005$$

$$\text{Cold class}(x) = P(\text{cold}) \times P(\text{congestion} = \text{True} | \text{cold})$$

$$\textcircled{3} \quad P(\text{Fever} = \text{False} \mid \text{cold}) \times P(\text{Chills} = \text{False} \mid \text{cold}) \times \\ P(\text{congestion} = \text{True} \mid \text{cold})$$

$$P(\text{cold} \mid X) = 2/6 \times 1/2 \times 1/2 \times 1/2 = 0.041$$

Hence the probability of cold class is greater than
probability of flu class ($0.041 > 0.005$)

\therefore Therefore patient P7 comes under cold class

b) Estimating probability / smoothing

$$P(X_i = a_i \mid c) = \frac{n_c^{x_i=a_i} + 1}{n_c + K}$$

To find

$$P(\text{congestion} = \text{True} \mid \text{cold}) = \frac{2+1}{2+2} = \frac{3}{4}$$

$$P(\text{congestion} = \text{False} \mid \text{cold}) = \frac{0+1}{2+2} = \frac{1}{4}$$

4) Logistic regression:-

Given $B_0 = -5.247$, $B_1 = 3.626$, also output

1 = No Bacteria, 0 = Bacteria present

$$P(Y=1|X) = \sigma(w \cdot x) = \frac{1}{1+e^{-(B_0+B_1x)}}$$

$$P(Y=0|X) = 1 - \sigma(w \cdot x) = \frac{e^{-(B_0+B_1x)}}{1+e^{-(B_0+B_1x)}}$$

To find

i) $P(\text{no bacteria} | \text{dosage} = 1.0)$

$$P(Y=0|X=1.0) \Rightarrow 0.165$$

ii) $P(\text{no bacterial dosage} = 2.5)$

$$P(Y=0|X=2.5) = \frac{1}{1+e^{-(5.247+3.626(2.5))}}$$

$$= P(Y=0|X=2.5) = 0.999$$

(iii) $P(\text{bacterial dosage } X = 2.0)$

$$P(Y=1|X=2.0) = \frac{e^{-(5.247+3.626(2))}}{1+e^{-(5.247+3.626(2))}}$$

$$P(Y=1|X=2.0) = 0.118$$

4) b) 80% percentile cutoff for dosage producing no bacteria (ie) to find the dosage (x) = ?

$$P(Y=1 \mid \text{dosage}(x)) = 0.8$$

no bacteria

$$P(Y=1 \mid \text{dosage}(x)) = \frac{1}{1+e^{-(B_0+B_1x)}} = 0.8$$

$$\left(\frac{1}{0.8}\right) - 1 = e^{-(B_0+B_1x)}$$

$$\ln(10 \cdot 0.25) = e^{-(B_0+B_1x)}$$

$$\ln(0.25) = \ln(e^{-(B_0+B_1x)})$$

$$\ln(0.25) = -(B_0+B_1x)$$

$$-\ln(0.25) \Rightarrow B_0+B_1x \quad \text{Multiplying by -1.}$$

$$\frac{-\ln(0.25) + 5.247}{3.626} = x$$

$$x = 1.828$$

Therefore $x = 1.828$ dosage, so 80th percentile
cutoff no bacteria

5) Given $x+y=1$
 $x+2y=2$ we got
 $x+3y=2$
 $x+4y=3$

$$\begin{bmatrix} A & X \\ \vdots & \vdots \\ 1 & 1 \\ 1 & 2 \\ 1 & 3 \\ 1 & 4 \end{bmatrix} \begin{bmatrix} x \\ y \end{bmatrix} = \begin{bmatrix} b \\ \vdots \\ 1 \\ 2 \\ 3 \\ 4 \end{bmatrix}$$

Using derivative formulae for OLS

$$x^* = (A^T A)^{-1} A^T b$$

Now to find $(A^T A)^{-1}$

$$(A^T A)^{-1} = \begin{bmatrix} 1 & 1 & 1 & 1 \\ 1 & 2 & 3 & 4 \end{bmatrix} \begin{bmatrix} 1 & 1 \\ 1 & 2 \\ 1 & 3 \\ 1 & 4 \end{bmatrix}^{-1}$$

$$= \begin{bmatrix} 1.5 & -0.5 \\ -0.5 & 0.2 \end{bmatrix} - \textcircled{1}$$

Now lets find $A^T b$

$$A^T b = \begin{bmatrix} 1 & 1 & 1 & 1 \\ 1 & 2 & 3 & 4 \end{bmatrix} \begin{bmatrix} 1 \\ 2 \\ 3 \\ 4 \end{bmatrix} = \begin{bmatrix} 8 \\ 23 \end{bmatrix} - \textcircled{2}$$

combine eq \textcircled{1} and \textcircled{2}

$$x^* = (A^T A)^{-1} A^T b$$

$$x^* = \begin{bmatrix} 0.5 \\ 0.6 \end{bmatrix}$$

$$y = 0.6 + 0.5x$$

(5)

Given $A = \begin{bmatrix} 0 & 1 \\ 1 & 3/2 \end{bmatrix}$

To find singular values of A , we should know λ_1, λ_2

also σ_1 and $\sigma_2 = ?$, $\sigma_1 = \sqrt{\lambda_1}$, $\sigma_2 = \sqrt{\lambda_2}$

Now let us find $A^T A$,

$$A^T A = \begin{bmatrix} 0 & 1 \\ 1 & 3/2 \end{bmatrix} \begin{bmatrix} 0 & 1 \\ 1 & 3/2 \end{bmatrix} = \begin{bmatrix} 1 & 3/2 \\ 3/2 & 13/4 \end{bmatrix}$$

To find λ_1 and $\lambda_2 = ?$

Let $\lambda^2 - S_1\lambda + S_2 = 0$ be the characteristic equations

$$\text{of } A^T A = 1 + 13/4 = 17/4$$

$S_2 = \text{determinant of matrix } (A^T A) (\det A = ad - bc)$

$$S_2 = 1 \times 13/4 - 9/4 = 4/4 = 1$$

Putting values of S_1 & S_2 in above equation

$$\lambda^2 - 17/4 + 1 = 0$$

From solving this equation we got value

$$\text{of } \lambda = 4, 0.25$$

$$\lambda_1 = 4, \lambda_2 = 1/4$$

also $\sigma_1 = \sqrt{\lambda_1}, \sigma = \sqrt{\lambda_2}$

$$\sigma_1 = \sqrt{4}, \sigma_2 = \sqrt{1/4}$$

We know that $S = \begin{bmatrix} \sigma_1 & 0 \\ 0 & \sigma_2 \end{bmatrix}$

$$\Sigma = \begin{bmatrix} \sqrt{4} & 0 \\ 0 & \sqrt{1/4} \end{bmatrix} \quad \Sigma = \begin{bmatrix} 2 & 0 \\ 0 & 0.5 \end{bmatrix}$$

Here the singular value of A is $\sigma_1 = \sqrt{4}$, $\sigma_2 = \sqrt{1/4}$
 Confirmation using matrix factorization multiplication

to prove

$$A = \begin{bmatrix} 0 & 1 \\ 1 & 3/2 \end{bmatrix} = \begin{bmatrix} 1/\sqrt{5} & -2/\sqrt{5} \\ 2/\sqrt{5} & 4/\sqrt{5} \end{bmatrix} \begin{bmatrix} 2 & 0 \\ 0 & 0.5 \end{bmatrix} \begin{bmatrix} 1/\sqrt{5} & 2/\sqrt{5} \\ 2/\sqrt{5} & -1/\sqrt{5} \end{bmatrix}$$

$$= \begin{bmatrix} 0.447 & -0.894 \\ 0.894 & 0.447 \end{bmatrix} \begin{bmatrix} 2 & 0 \\ 0 & 0.5 \end{bmatrix} \begin{bmatrix} 0.447 & 0.894 \\ 0.894 & -0.447 \end{bmatrix}$$

$$= \begin{bmatrix} 0 & 1 \\ 1 & 1.5 \end{bmatrix}$$

Hence $A = U\Sigma V^T$ proved.

6) ii) Outer product form of SVD, by taking the rank-1 approximation for matrix A.

Rank k=1 [We only consider largest singular value $\sigma_1 = \sqrt{4}$]

$$A_k = \sigma_1 U_1 V_1^T$$

$$A_k = \sqrt{4} \begin{bmatrix} 1/\sqrt{5} \\ 2/\sqrt{5} \end{bmatrix} \begin{bmatrix} 1/\sqrt{5} & 2/\sqrt{5} \end{bmatrix}$$

(6)

$$A = \begin{bmatrix} 0.4 & 0.8 \\ 0.8 & 1.6 \end{bmatrix}$$

6) iii) To find SVD for A^{-1}
 Formulae $A^{-1} = V \Sigma^{-1} U^T$

$$A^{-1} = \begin{bmatrix} -3/2 & 1 \\ 1 & 0 \end{bmatrix} \quad U^T = \begin{bmatrix} 1/\sqrt{5} & 2/\sqrt{5} \\ -2/\sqrt{5} & 1/\sqrt{5} \end{bmatrix}$$

$$V = \begin{bmatrix} 1/\sqrt{5} & 2/\sqrt{5} \\ 2/\sqrt{5} & -1/\sqrt{5} \end{bmatrix} \quad \Sigma^{-1} = \begin{bmatrix} 1/4 & 0 \\ 0 & \sqrt{4} \end{bmatrix}$$

Putting in formulae $A^{-1} = V \Sigma^{-1} U^T$

$$A^{-1} = \begin{bmatrix} -3/2 & 1 \\ 1 & 0 \end{bmatrix} = \begin{bmatrix} 1/\sqrt{5} & 2/\sqrt{5} \\ 2/\sqrt{5} & -1/\sqrt{5} \end{bmatrix} \begin{bmatrix} \sqrt{4} & 0 \\ 0 & \sqrt{4} \end{bmatrix} \begin{bmatrix} 1/\sqrt{5} & 2/\sqrt{5} \\ -2/\sqrt{5} & 1/\sqrt{5} \end{bmatrix}$$

$$A^{-1} = \begin{bmatrix} -1.5 & 1 \\ 1 & 0 \end{bmatrix} = \begin{bmatrix} 0.223 & 1.788 \\ 0.447 & -0.894 \end{bmatrix} \begin{bmatrix} 0.447 & 0.894 \\ -0.894 & 0.447 \end{bmatrix}$$

$$A^{-1} = \begin{bmatrix} -1.5 & 1 \\ 1 & 0 \end{bmatrix} = \begin{bmatrix} -1.49 & 0.99 \\ 0.99 & 0 \end{bmatrix} \cong \begin{bmatrix} -1.5 & 1 \\ 1 & 0 \end{bmatrix}$$

Hence proved

$$A^T = U \Sigma^{-1} V$$

(6)

7) A) Given 2 cluster m_1, m_2 with center points

$$m_1 = (1, 1, 1) \text{ and } m_2 = (1, 2, 0)$$

Data points	distance		Cluster
	$m_1 = (1, 1, 1)$	$m_2 = (1, 2, 0)$	
$x_1 = (1, 1, 1)$	0	$\sqrt{2}$	m_1
$x_2 = (1, 1, 4)$	3	4.12	m_2
$x_3 = (1, 1, 2)$	4.12	4.58	m_2 m_1
$x_4 = (1, 2, 0)$	$\sqrt{2}$	0	m_2

Formulae to find distance b/w 2 points for m_1 :

$$d(p_1, p_2) = \sqrt{(m_{1,1} - x_{1,1})^2 + (m_{1,2} - x_{1,2})^2 + (m_{1,3} - x_{1,3})^2}$$

Similarly for m_2

$$d(p_1, p_2) = \sqrt{(m_{2,1} - x_{1,1})^2 + (m_{2,2} - x_{1,2})^2 + (m_{2,3} - x_{1,3})^2}$$

Let us find the distance b/w 2 points for

from m_1 (center) to all x points

$$= \sqrt{(1-1)^2 + (1-1)^2 + (1-1)^2} = 0$$

similarly distance from m_2 to all x_1 points

$$= \sqrt{(1-1)^2 + (2-1)^2 + (0-1)^2} = \sqrt{2}$$

$$\text{for } x_4: \sqrt{(1-1)^2 + (1-2)^2 + (0-0)^2} = \sqrt{2}$$

$$\text{and } \sqrt{(1-1)^2 + (2-2)^2 + (0-0)^2} = 0$$

Q) Since we got x_1, x_2, x_3 comes under m_1 cluster & x_4
 is m_2 cluster
 So we have to find new centroid for $m_1 = ?$ & for
 m_2 the centroid will be same (since only 1
 instance (x_4) comes under it)
 New centroid for $m_1 = [2/3, 1, 2/3]$
 ie $m_1 \approx (2.3, 1, 2.33)$ remains same

$$m_2 \approx (1, 2, 0)$$

$$\text{7) b) mean entropy } (C) = \sum_{i=1}^k \frac{n_i}{n} \text{Entropy}(C_i)$$

$$\text{where } (C_i) = \sum_{j=1}^n p_{ij} \log_2 p_{ij}$$

Now to find entropy C_1 & C_2

$$\begin{aligned} \text{Entropy } C_1 &= (2/3 \log_2 2/3 + 1/3 \log_2 1/3) \\ &= -(-0.929) \end{aligned}$$

$$\text{Entropy}(C_1) = 0.929, \text{entropy}(C_2) = -(\log_2 1) = 0$$

$$\begin{aligned} \text{Find mean entropy } (C) &= \frac{3}{4} \text{Entropy}(C_1) + \frac{1}{4} \text{Entropy}(C_2) \\ &= \frac{3}{4} * 0.929 + 0 \end{aligned}$$

$$\text{Entropy}(C) = 0.696$$

Part 2 Programming assignment

Assignment 1 :

Given us

The initial starting points for the K cluster means can be K randomly selected data points. You should have an option to run the algorithm r times from r different randomly chosen initializations

```
LOWEST SUM OF SQUARE ERROR FOR RUN [10]: 946.6381021673415
SUM OF SQUARE ERROR RUN [10]: 1046.1569429493627
```

```
PLOTTING GRAPH FOR LOWEST SUM OF SQUARE ERROR: 946.6381021673415
```

The plot is included in a separate file called K-Means with screenshots pasted

Assignment 2 :

The main distinction between K-means and C-means is that in C-means, a data point can belong to multiple clusters. The user provides the number of clusters and interactions as input. Random membership weights are assigned to each data point. The membership grade and centroid for each cluster are computed using the fuzzifier parameter $m = 2$. The membership grade quantifies the degree to which an individual data point can belong to a cluster.

```
Epsilon reached, stopping early
New best square error: 2
All squareErrors sums:
[1196.8595182099502, 1177.3585820963945]
Best squareError sum:
1177.3585820963945
```

The plot is included in a separate file called C-Means with screenshots pasted

