

1. Problem 1: Probability and Expectation (22 points)

- In a bag there are 50 balls, each having one of three colors: red (25 balls), blue (15 balls), and green (10 balls). When a ball is drawn from the bag, you get different points for every color:
 - Red balls are worth 3 points each.
 - Green balls are worth 5 points each.
 - Blue balls are worth 10 points each.

1. Suppose that you draw a ball at random from the bag. What is the expected total point value of the ball you are drawing?

$$P(R) = 25/50 = 1/2 * 3 = 1.5$$

$$P(G) = 10/50 = 1/5 * 5 = 1$$

$$P(B) = 15/50 = 3/10 * 10 = 3$$

$$E(x) = 1.5 + 1 + 3 = 5.5. \quad E(x) = 5.5$$

2. I. X is the random variable that denotes the number of points you get by drawing a ball at random, what is the variance of X ?

$$\begin{aligned} E(X^2) &= (3^2 * 1/2) + (5^2 * 1/5) + (10^2 * 3/10) \\ &= 4.5 + 5 + 30 = 39.5 \end{aligned}$$

$$\begin{aligned} \text{Var}(X) &= E[X^2] - E^2[X] \\ &= 39.5 - (5.5)^2 = 9.25 \quad \text{Var}(X) = 9.25 \end{aligned}$$

3. Three balls were picked at random from the bag and it turned out to be two red balls and one green ball. Now, you are drawing another ball from the remaining ones. What is the expected total point value of the ball you are drawing?

$$P(R) = 23/47 * 3 \approx 1.468$$

$$P(G) = 9/47 * 5 \approx 0.957$$

$$P(B) = 15/47 * 10 \approx 3.191.$$

$$E(x) = 1.468 + 0.957 + 3.191 \approx 5.616 \quad E(x) \approx 5.616$$

4. Suppose that X is a random variable, with expectation $E(X) = 10$ and variance $\text{Var}(X) = 2$. What is the value of $E[X(X - 1)]$?

$$E[X(X - 1)] = E(X^2) - E(X)$$

$$\text{Var}(X) = E[X^2] - E^2[X]$$

$$2 = E[X^2] - 10^2 \quad E[X^2] = 102$$

$$E[X(X - 1)] = 102 - 10 = 92 \quad E[X(X - 1)] = 92$$

2. Problem 2: Conditional Probability and Bayes Theorem (18 points)

Probability theory and the Bayes theorem serve as the basis of important machine learning techniques such as linear discriminant analysis (LDA), density estimation, and Naive Bayes (NB) classifiers, which are very effective in text classification settings.

A spam filter is designed by looking at commonly occurring phrases in spam. Suppose that 80% of email is spam. In 10% of the spam emails, the phrase "free money" is used, whereas this phrase is only used in 1% of non-spam emails.

1. What is the probability that an email has the phrase "free money"?

$$80/100 = P(S)$$

$$20/100 = P(N)$$

$$0.01 = P(F | N)$$

$$0.1 = P(F | S)$$

$$P(F) = (0.1 * 0.8) + (0.01 * 0.2) = 0.08 + 0.002 \quad P(F) = 0.082 = 8.2\%$$

2. A new email has just arrived, which does mention "free money". What is the probability that it is spam?

$$P(S | F) = \frac{P(F|S)P(S)}{P(F)}$$

$$= \frac{0.1 * 0.8}{0.082} = 0.9756097561 \quad P(S | F) \approx 0.976 = 97.6\%$$

3. Another email does not have the phrase "free money". What is the probability that it is not spam?

$$P("Not Spam" | "Not Free Money") = \frac{P(\neg F|N)P(N)}{P(\neg F)}$$

$$= \frac{0.99 * 0.2}{0.918}$$

$$= 0.2156862745 \quad P(N | \neg F) \approx 0.216 = 21.6\%$$

3. Problem 3: Matrices and vectors (20 points)

Since data is represented in matrices, many of the techniques that we will see, from linear regression to dimensionality reduction and neural networks, will boil down to matrix operations. Optimized matrix operations from modern packages allow us to train and test algorithms much faster than in an iterative approach. Define the following two matrices of size 4×4 :

$$D_1 = \begin{bmatrix} 1 & 0 & 2 & 0 \\ 0 & 2 & 4 & -2 \\ -1 & 0 & -3 & 1 \\ 0 & -1 & -1 & 0 \end{bmatrix}$$

$$D_2 = \begin{bmatrix} 1 & -1 & 0 & 2 \\ -1 & 0 & 1 & -1 \\ 0 & -2 & 0 & 1 \\ 3 & 0 & -1 & 0 \end{bmatrix}$$

1. Are the 4 columns in each of the matrix linearly independent or linearly dependent?

D1 = Linearly dependent

D2 = Linearly Independent

1)

$$D_1 = \begin{bmatrix} 1 & 0 & 2 & 0 \\ 0 & 2 & 4 & -2 \\ -1 & 0 & -3 & 1 \\ 0 & -1 & -1 & 0 \end{bmatrix} \quad \text{expand by col 1}$$

$$1 \cdot \begin{bmatrix} 2 & 4 & -2 \\ 0 & -3 & 1 \\ -1 & -1 & 0 \end{bmatrix} - 1 \det \begin{bmatrix} 0 & 2 & 0 \\ 2 & 4 & -2 \\ -1 & -1 & 0 \end{bmatrix}$$

$$\left[2 \det \begin{bmatrix} -3 & 1 \\ -1 & 0 \end{bmatrix} + (-1) \det \begin{bmatrix} 4 & -2 \\ -3 & 1 \end{bmatrix} \right] - \left[2 \det \begin{bmatrix} 2 & -2 \\ -1 & 0 \end{bmatrix} \right]$$

$$2(0+1) + -1(4-6) \quad 2(0-2)$$

$$2+2=4 \quad -1(-4)=4$$

$$4-4=\underline{\underline{0}} \quad \boxed{\det(D_1)=0 \therefore \text{linearly dependent}}$$

$$D_2 = \begin{bmatrix} 1 & -1 & 0 & 2 \\ -1 & 0 & 1 & -1 \\ 0 & -2 & 0 & 1 \\ 3 & 0 & -1 & 0 \end{bmatrix} \quad \text{expand by col 2}$$

$$(-1) \det \begin{bmatrix} -1 & 1 & -1 \\ 0 & 0 & 1 \\ 3 & -1 & 0 \end{bmatrix} - (2) \det \begin{bmatrix} 1 & 0 & 2 \\ -1 & 1 & -1 \\ 3 & -1 & 0 \end{bmatrix}$$

$$\left[1 \begin{bmatrix} -1 & 1 \\ 3 & -1 \end{bmatrix} \right] - 2 \left[\det \begin{bmatrix} 1 & 2 \\ 3 & 0 \end{bmatrix} - 1 \det \begin{bmatrix} 1 & 2 \\ -1 & -1 \end{bmatrix} \right]$$

$$-1(1-3) \quad (0-6) - (-1)(1+2)$$

$$2 \quad -6+1$$

$$2-(-5)=\underline{\underline{10}} \quad \boxed{\det(D_2) \text{ Linearly INDEPENDENT}}$$

2. For each of the matrices above, find the maximum number of columns that are linearly independent.

$\det(D_1) = 0$ and linearly dependent, therefore rank < 4 ,

Since there exists a nonzero 3×3 subset, **the maximum number of columns that are linearly independent is 3**

$$\cdot \det(D_2) = -8 \neq 0$$

Nonzero determinant = **4 linearly independent columns**

3. What is the rank of each matrix? Use the procedure we discussed in class to compute the rank of each matrix, and show the steps you performed.

$$\text{Row 1: } 1 \cdot c_1 + 0 \cdot c_2 + 2 \cdot c_3 + 0 \cdot c_4 = 0 \implies c_1 + 2c_3 = 0$$

$$\text{Row 2: } 0 \cdot c_1 + 2 \cdot c_2 + 4 \cdot c_3 - 2 \cdot c_4 = 0 \implies 2c_2 + 4c_3 - 2c_4 = 0$$

$$\text{Row 3: } -1 \cdot c_1 + 0 \cdot c_2 - 3 \cdot c_3 + 1 \cdot c_4 = 0 \implies -c_1 - 3c_3 + c_4 = 0$$

$$\text{Row 4: } 0 \cdot c_1 - 1 \cdot c_2 - 1 \cdot c_3 + 0 \cdot c_4 = 0 \implies -c_2 - c_3 = 0$$

$$c_1 + 2c_3 \mid 2c_2 + 4c_3 - 2c_4 \mid -c_1 - 3c_3 + c_4 \mid -c_2 - c_3$$

$$\text{Row 1: } c_1 = -2c_3$$

$$\text{Row 4: } c_2 = -c_3 \quad \text{Row 2: } 0 = -2(c_3) + 4c_3 - 2c_4 \Rightarrow 2c_3 - 2c_4 = 0 \Rightarrow \underline{c_3 = c_4}$$

$$\text{Row 3: } -c_1 - 3c_3 + c_4 = 0 \Rightarrow 2c_3 - 3c_3 + c_3 = 0 \Rightarrow 3c_3 - 3c_3 = 0 \Rightarrow 0 = 0$$

$$c_3 = x$$

$$c_4 = x$$

$$c_1 = -2x, c_2 = -x$$

Since there is 1 free unknown, 4 columns - 1 free = 3, therefore, **Rank(D1) = 3**

$$\text{Row 1: } (c_1 - c_2 + 2c_4 = 0)$$

$$\text{Row 2: } (-c_1 + c_3 - c_4 = 0)$$

$$\text{Row 3: } (-2c_2 + c_4 = 0)$$

$$\text{Row 4: } (3c_1 - c_3 = 0)$$

$$\text{Row 4: } c_3 = 3c_1 \quad \text{Row 3: } c_4 = 2c_2$$

$$\text{Row 2: } -c_1 + c_3 - c_4 = 0 \implies -c_1 + 3c_1 - 2c_2 = 0 \implies 2c_1 - 2c_2 = 0 \implies c_2 = c_1$$

$$\text{Row 1: } c_1 - c_2 + 2c_4 = c_1 - c_1 + 2(2c_2) = 4c_1 = 0 \implies c_1 = 0$$

$$c_1 = 0 \implies c_2 = 0, \quad c_3 = 0, \quad c_4 = 0$$

The solution is trivial, this means there are no free unknown,

therefore **Rank(D2) = 4 - 0 = 4**

4. Assume that we define a column vector θ . Which of the products θD_1 , $D_1\theta$, $\theta^T D_1$, and $D_1\theta^T$ can be computed? For the ones that can be computed, please write the result and its dimension. Please perform these computations by hand and show your work. Do not use a calculator for this problem.

$$\theta = \begin{bmatrix} -1 \\ 0 \\ 1 \\ 2 \end{bmatrix}$$

4.

$$\theta = (4 \times 1) \quad D_1 = (4 \times 4)$$

a) $\theta D_1 = (4 \times 1)(4 \times 4)$
 does NOT match
 Cannot be computed

b) $D_1 \theta = (4 \times 4)(4 \times 1)$
 matches ✓ can be computed

$$\begin{bmatrix} 1 & 0 & 2 & 0 \\ 0 & 2 & 4 & -2 \\ -1 & 0 & -3 & 1 \\ 0 & -1 & -1 & 0 \end{bmatrix} \times \begin{bmatrix} -1 \\ 0 \\ 1 \\ 2 \end{bmatrix} = \begin{bmatrix} -1 + 0 + 2 + 0 \\ 0 + 0 + 4 - 4 \\ 1 + 0 - 3 + 2 \\ 0 + 0 - 1 + 0 \end{bmatrix} = \begin{bmatrix} 1 \\ 0 \\ 0 \\ -1 \end{bmatrix}$$

c) $\theta^T D_1 = (1 \times 4)(4 \times 4)$
 matches ✓ can be computed

$$\begin{bmatrix} -1 & 0 & 1 & 2 \end{bmatrix} \times \begin{bmatrix} 1 & 0 & 2 & 0 \\ 0 & 2 & 4 & -2 \\ -1 & 0 & -3 & 1 \\ 0 & -1 & -1 & 0 \end{bmatrix} = \begin{bmatrix} (-1 + 0 - 1 + 0), \\ (0 + 0 + 0 - 2), \\ (-2 + 0 - 3 - 2), \\ (0 + 0 + 1 + 0) \end{bmatrix} = \underline{\begin{bmatrix} -2 & -2 & -7 & 1 \end{bmatrix}}$$

d) $D_1 \theta^T = (4 \times 4)(1 \times 4)$
 does NOT Match
 Cannot be computed

4. Problem 4: Matrix transpose and inverse (20 points)

1. Anton speaks French and German; Geraldine speaks English, French and Italian; James speaks English, Italian, and Spanish; Lauren speaks all the languages the others speak except French; and no one speaks any other language. Make a matrix $A = \{a_{ij}\}$ with rows representing the four people mentioned and columns representing the languages they speak. Put $a_{ij} = 1$ if person i speaks language j and $a_{ij} = 0$ otherwise.

Compute the matrices AA^T and $A^T A$. Explain the significance of the matrices AA^T and $A^T A$. What does each entry in these two matrices represent?

Anton = F, G

Geraldine = E, F, I

James = E, I, S

Lauren = G, E, I, S

$A = \{a_{ij}\}$

$a_{ij} = 1 \rightarrow \text{Person } i \text{ speaks language } j \text{ and } a_{ij} = 0$

$AA^T:$

$$A = \begin{bmatrix} F & G & E & I & S \\ 1 & 1 & 0 & 0 & 0 \\ 1 & 0 & 1 & 1 & 0 \\ 0 & 0 & 1 & 1 & 1 \\ 0 & 1 & 1 & 1 & 1 \end{bmatrix} \cdot A^T = \begin{bmatrix} 1 & 1 & 0 & 0 \\ 1 & 0 & 0 & 0 \\ 0 & 1 & 1 & 1 \\ 0 & 1 & 1 & 1 \\ 0 & 0 & 1 & 1 \end{bmatrix} = \begin{bmatrix} 1+1, 1, 0, 1 \\ 1, 1+2, 2, 2 \\ 0, 2, 3, 3 \\ 1, 2, 3, 4 \end{bmatrix}$$

$(4 \times 6) \times (5 \times 4)$

A	G	J	L
2	1	0	1
1	3	2	2
0	2	3	3
1	2	3	4

This is a person x person chart. It shows the number of languages 2 people speak in common

$A^TA:$

$$A^T = \begin{bmatrix} 1 & 1 & 0 & 0 \\ 1 & 0 & 0 & 1 \\ 0 & 1 & 1 & 1 \\ 0 & 1 & 1 & 1 \\ 0 & 0 & 1 & 1 \end{bmatrix} \times A = \begin{bmatrix} F & G & E & I & S \\ 1 & 1 & 0 & 0 & 0 \\ 1 & 0 & 1 & 1 & 0 \\ 0 & 1 & 1 & 1 & 1 \\ 0 & 0 & 1 & 1 & 1 \\ 0 & 1 & 1 & 1 & 1 \end{bmatrix} = \begin{bmatrix} 2, 1, 1, 1, 1, 0 \\ 1, 1+1, 1, 1, 1, 2 \\ 1, 1, 3, 3, 3, 2 \\ 1, 1, 3, 3, 3, 2 \\ 0, 1, 2, 2, 2, 2 \end{bmatrix}$$

$(5 \times 4) \times (4 \times 6)$

F	G	E	I	S
2	1	1	1	0
1	2	1	1	1
1	1	3	3	2
1	1	3	3	2
0	1	2	2	2

This is a language x language chart. It shows how many people speak both 2 languages.

2. Generate at random 3 matrices of size 3×3 and fill each entry with a random integer chosen from -10 to 10. Use an existing package to compute the inverse of each matrix if it exists. Include in your report:

- Each of the 3 matrices
- The inverses computed with the package
- Compute the product of each matrix with its inverse to check that you obtain the identity matrix.

Answer:

<https://github.com/haderie/DS4400/blob/main/HW1/Problem4.ipynb>

5. Problem 5: Average, variance, and correlation (20 points)

The dataset for this assignment is available [here](#). The prediction task is to predict the price of a house (column **price**) given the other features. Please ignore the columns **id** and **date**, as well as the categorical column **zipcode**. File **kc_house_data.csv** includes all the records in the dataset. You should use the entire dataset for the assignment. You can also find a Word document including the feature description in the same folder. In this problem, we will perform some exploratory data analysis using the house price dataset.

1. For each feature, write code to compute the average value, the min and max values, as well as its variance.

Which features have the lowest and the highest average? Include the feature name and their average values for the features with the lowest and highest average.

Which features have the lowest and the highest variance? Include the feature name and their variance values for the features with the lowest and highest variance.

2. Compute the correlation coefficient of each feature with the response. Include a table with the correlation coefficient of each feature with the response. Which features are positively correlated (i.e., have positive correlation coefficient) with the response?
3. Which feature has the highest positive correlation with the response? 3. Were you able to find any features with a negative correlation coefficient with the response?

Answer:

<https://github.com/haderie/DS4400/blob/main/HW1/Problem5.ipynb>