

IE30301-Datamining Assignment 3 (70 Points)

Eldor Fozilov

April 27, 2021

Exercise 1

Write a detailed description of the following concepts and their differences. (Explain at least 2 lines about each concepts and write differences in 1 sentence. If not, there are 3 points deduction per problem. [20 pts, 5 pts for each.]

1. Likelihood & Probability

Probability, in simple words, can be defined as the chance that a given event will occur. It is calculated as the ratio of the number of outcomes in a sample space representing a given event of interest to the total number of possible outcomes in that sample space. The word **likelihood** is often interchangeably used with the word **probability** in daily conversations. However, it has a different meaning in statistics and machine learning. In those areas, **likelihood** generally concerned with finding the best distribution of the data given a sample data set. It is represented as a function (called **likelihood function**), which calculates the joint probability distribution of given data as a function of model parameters. That function can be thought as the probability of obtaining the given data given the parameters.

2. Feature Selection & Feature Extraction

Feature selection is referred to the selection of a subset of features out of the original features with the aim of simplifying a model, reducing computational cost and the influence of noise, caused by redundant features in the model. So, the result of successful feature selection can be an improvement in the model accuracy and speed.

Feature extraction serves similar purposes mentioned above, but it tries to achieve it in a different way compared to **feature selection**: instead of choosing which variables to include in a model, feature extraction methods create new variables from all given feature variables by combining them in various ways so that those new variables retain most of the useful information. For example PCA creates new variables by taking a **linear** combination of given feature variables.

3. PC Score & PC loading

We know the fact that **PCA** produces a low-dimensional representation of a data set by finding a sequence of linear combinations of the variables that retain maximum variance, and are uncorrelated with each other. The particular values in those derived variable vectors are called **PC scores**.

To find those **PC scores**, we need to take a dot product between each normalized sample

record, which can be thought as being a vector, and a parameter vector. The values in that parameter vector are called **PC loadings**, which represent the contribution of each feature variable to the computation of **PC score**. Those **PC loadings** can also be interpreted as correlations between feature variables and the loading vector (principal axis) they belong to. So we can see that PC score and PC loading are closely related, but not the same thing.

4. Newton-Raphson Method & Gradient Descent

Gradient descent is an iterative optimization algorithm for finding a local minimum of a differentiable function using its first-order derivative. It is grounded in the idea that if a derivative of a multivariate function f exists at a point \mathbf{x} , then the direction of fastest decrease/descent in the function value will happen when if one goes from the point \mathbf{x} in the direction of the negative gradient of f at \mathbf{x} , $-\nabla f(\mathbf{x})$. That negative gradient can be multiplied by a some value to increase or decrease the step size (change in value of \mathbf{x}).

Newton-Raphson method is an iterative method used for finding the roots of a twice differentiable function f . The main difference between those two methods is that gradient descent algorithm is parametric in nature (it usually requires a learning rate μ), while the other method does not need such a parameter, and thus we can apply it without worrying for changing any hyperparameter.

Exercise 2

In multiple linear regression, we can estimate $\hat{\beta}$ as follows by least square method.

$$\hat{\beta} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} \in \mathbb{R}^{(p+1) \times 1}, \quad \mathbf{X} \in \mathbb{R}^{N \times (p+1)}, \quad \mathbf{y} \in \mathbb{R}^{N \times 1} \quad (2.1)$$

The regression model can be derived as follows:

$$\hat{\mathbf{y}} = \mathbf{X} \hat{\beta} = \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} = \mathbf{H} \mathbf{y} \quad (2.2)$$

And in the above equation, $\mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T$ is specifically referred to as \mathbf{H} , hat matrix.

$$\mathbf{H} = \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \quad (2.3)$$

2.1

Show that \mathbf{H} is symmetric ($\mathbf{H}^T = \mathbf{H}$) and idempotent ($\mathbf{H}^2 = \mathbf{H}$). [3 pts]

Answer:

First we will show that the matrix \mathbf{H} is symmetric.

$$\mathbf{H}^T = [\mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T]^T$$

We will use the following property of the transpose applied when multiplying two matrices:

$$(\mathbf{AB})^T = \mathbf{B}^T \mathbf{A}^T \quad (2.4)$$

Let's label the matrix $(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T$ as M. Then we write \mathbf{H}^T as

$$\mathbf{H}^T = [\mathbf{X} \mathbf{M}]^T = \mathbf{M}^T \mathbf{X}^T = [(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T]^T \mathbf{X}^T = [(\mathbf{X}^T)^T ((\mathbf{X}^T \mathbf{X})^{-1})^T] \mathbf{X}^T$$

We know that $(\mathbf{X}^T)^T = \mathbf{X}$. We can also show that $((\mathbf{X}^T \mathbf{X})^{-1})^T = (\mathbf{X}^T \mathbf{X})^{-1}$.

Let $\mathbf{A} = (\mathbf{X}^T \mathbf{X})^{-1}$. Then $\mathbf{A}^{-1} \mathbf{A} = \mathbf{I}$ and $\mathbf{A} \mathbf{A}^{-1} = \mathbf{I}$. If we take the transpose of both sides in those equations and apply (2.4), then we get the following:

$$\mathbf{A}^T (\mathbf{A}^{-1})^T = \mathbf{I}^T = \mathbf{I}, \quad (\mathbf{A}^{-1})^T \mathbf{A}^T = \mathbf{I}^T = \mathbf{I}$$

From the above equations, we can see that $(\mathbf{A}^T)^{-1} = (\mathbf{A}^{-1})^T$.

$$\text{So, } ((\mathbf{X}^T \mathbf{X})^{-1})^T = ((\mathbf{X}^T \mathbf{X})^T)^{-1} = (\mathbf{X}^T \mathbf{X})^{-1}$$

Thus, \mathbf{H}^T can be expresses as

$$\mathbf{H}^T = [(\mathbf{X}^T)^T ((\mathbf{X}^T \mathbf{X})^{-1})^T] \mathbf{X}^T = \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T = \mathbf{H}$$

Now, we will show that the matrix \mathbf{H} is idempotent.

$$\begin{aligned} \mathbf{H}^2 &= \mathbf{H} \mathbf{H} = (\mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T) (\mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T) = \\ &= \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} (\mathbf{X}^T \mathbf{X}) (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T = \mathbf{X} \mathbf{I} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T = \\ &= \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T = \mathbf{H} \end{aligned}$$

2.2

An estimator of a given parameter is said to be unbiased if its expected value is equal to the true value of the parameter. ($E[\hat{\theta}] = \theta$)

Show that $MSE = \hat{\sigma}^2 = \frac{SSE}{N - p - 1} = \frac{(\mathbf{y} - \hat{\mathbf{y}})^T (\mathbf{y} - \hat{\mathbf{y}})}{N - p - 1}$ **is unbiased estimator** through using following properties (2.5) - (2.9) & results of above problem 2.1 [10 pts]

If c is a scalar, and \mathbf{A} is a $n \times n$ square matrix, ($c \in \mathbb{R}, \mathbf{A} \in \mathbb{R}^{n \times n}$)

$$c = \text{Tr}(c), \quad \text{Tr}(c\mathbf{A}) = c\text{Tr}(\mathbf{A}) \quad (2.5)$$

If \mathbf{A}, \mathbf{B} are $n \times n$ square matrix that have same dimension, ($\mathbf{A}, \mathbf{B} \in \mathbb{R}^{n \times n}$)

$$\text{Tr}(\mathbf{A} + \mathbf{B}) = \text{Tr}(\mathbf{A}) + \text{Tr}(\mathbf{B}), \quad \text{Tr}(\mathbf{A} - \mathbf{B}) = \text{Tr}(\mathbf{A}) - \text{Tr}(\mathbf{B}) \quad (2.6)$$

If \mathbf{A} is a $n \times n$ square matrix, $\mathbf{A} \in \mathbb{R}^{n \times n}$

$$E[\text{Tr}(\mathbf{A})] = \text{Tr}(E[\mathbf{A}]) \quad (2.7)$$

If \mathbf{A} is a $n \times m$ matrix and \mathbf{B} is a $m \times n$ matrix ($\mathbf{A} \in \mathbb{R}^{n \times m}, \mathbf{B} \in \mathbb{R}^{m \times n}$)

$$\text{Tr}(\mathbf{AB}) = \text{Tr}(\mathbf{BA}) \quad (2.8)$$

If \mathbf{x} is random vector,

$$\text{Var}[\mathbf{x}] = E[\mathbf{xx}^T] - E[\mathbf{x}]E[\mathbf{x}]^T, \quad E[\mathbf{xx}^T] = \text{Var}[\mathbf{x}] + E[\mathbf{x}]E[\mathbf{x}]^T \quad (2.9)$$

Answer:

Finding the expectation of **SSE** will give us clue to the expectation of **MSE** since **MSE** is just a SSE divided by $N - p - 1$. So we will tackle it first.

$$\begin{aligned} E[\text{SSE}] &= E[(\mathbf{y} - \hat{\mathbf{y}})^T(\mathbf{y} - \hat{\mathbf{y}})] = E[(\mathbf{y} - \mathbf{Hy})^T(\mathbf{y} - \mathbf{Hy})], \quad \text{where } \mathbf{H} = \mathbf{X}(\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T \\ E[(\mathbf{y} - \mathbf{Hy})^T(\mathbf{y} - \mathbf{Hy})] &= E[((\mathbf{I} - \mathbf{H})\mathbf{y})^T((\mathbf{I} - \mathbf{H})\mathbf{y})] = \\ &= E[\text{Tr}((\mathbf{I} - \mathbf{H})\mathbf{y})^T((\mathbf{I} - \mathbf{H})\mathbf{y})] \quad (\because \text{property (2.5)}) \\ E[\text{Tr}((\mathbf{I} - \mathbf{H})\mathbf{y})^T((\mathbf{I} - \mathbf{H})\mathbf{y})] &= E[\text{Tr}(\mathbf{y}^T(\mathbf{I} - \mathbf{H})^T(\mathbf{I} - \mathbf{H})\mathbf{y})] = \\ &= E[\text{Tr}(\mathbf{y}^T(\mathbf{I} - \mathbf{H})(\mathbf{I} - \mathbf{H})\mathbf{y})] \quad (\because \mathbf{H} \text{ and } \mathbf{I} \text{ are symmetric and thus } (\mathbf{I} - \mathbf{H}) \text{ is also symmetric}) \\ E[\text{Tr}(\mathbf{y}^T(\mathbf{I} - \mathbf{H})(\mathbf{I} - \mathbf{H})\mathbf{y})] &= E[\text{Tr}(\mathbf{y}^T(\mathbf{I} - \mathbf{H})^2\mathbf{y})] = \\ &= E[\text{Tr}(\mathbf{y}^T(\mathbf{I} - \mathbf{H})\mathbf{y})] \quad (\because \mathbf{I} \text{ and } \mathbf{H} \text{ are idempotent and thus } (\mathbf{I} - \mathbf{H}) \text{ is idempotent}) \\ E[\text{Tr}(\mathbf{y}^T(\mathbf{I} - \mathbf{H})\mathbf{y})] &= E[\text{Tr}((\mathbf{I} - \mathbf{H})\mathbf{y}\mathbf{y}^T)] \quad (\because \text{property (2.8)}) \\ E[\text{Tr}((\mathbf{I} - \mathbf{H})\mathbf{y}\mathbf{y}^T)] &= \text{Tr}(E[(\mathbf{I} - \mathbf{H})\mathbf{y}\mathbf{y}^T]) \quad (\because \text{property (2.7)}) \\ \text{Tr}(E[(\mathbf{I} - \mathbf{H})\mathbf{y}\mathbf{y}^T]) &= \text{Tr}((\mathbf{I} - \mathbf{H})E[\mathbf{y}\mathbf{y}^T]) \quad (\because \mathbf{X} \text{ is considered to be given, thus } \mathbf{H} \text{ is a constant}) \\ \text{Tr}((\mathbf{I} - \mathbf{H})E[\mathbf{y}\mathbf{y}^T]) &= \text{Tr}((\mathbf{I} - \mathbf{H})(\text{Var}[\mathbf{y}] + E[\mathbf{y}]E[\mathbf{y}]^T)) \quad (\because \text{property (2.9)}) \\ \text{Tr}((\mathbf{I} - \mathbf{H})(\text{Var}[\mathbf{y}] + E[\mathbf{y}]E[\mathbf{y}]^T)) &= \text{Tr}((\mathbf{I} - \mathbf{H})(\sigma^2\mathbf{I} + E[\mathbf{y}]E[\mathbf{y}]^T)) = \\ &= \text{Tr}((\mathbf{I} - \mathbf{H})\sigma^2\mathbf{I} + (\mathbf{I} - \mathbf{H})E[\mathbf{y}]E[\mathbf{y}]^T) = \text{Tr}((\mathbf{I} - \mathbf{H})\sigma^2\mathbf{I}) + \text{Tr}((\mathbf{I} - \mathbf{H})E[\mathbf{y}]E[\mathbf{y}]^T) \quad (\because \text{property (2.6)}) \\ \text{Tr}((\mathbf{I} - \mathbf{H})\sigma^2\mathbf{I}) + \text{Tr}((\mathbf{I} - \mathbf{H})E[\mathbf{y}]E[\mathbf{y}]^T) &= \text{Tr}(\sigma^2\mathbf{I} - \sigma^2\mathbf{H}) + \text{Tr}((\mathbf{I} - \mathbf{H})E[\mathbf{y}]E[\mathbf{y}]^T) = \\ &= \text{Tr}(\sigma^2\mathbf{I}) - \text{Tr}(\sigma^2\mathbf{H}) + \text{Tr}((\mathbf{I} - \mathbf{H})E[\mathbf{y}]E[\mathbf{y}]^T) \quad (\because \text{property (2.6)}) \\ \text{Tr}(\sigma^2\mathbf{I}) - \text{Tr}(\sigma^2\mathbf{H}) + \text{Tr}((\mathbf{I} - \mathbf{H})E[\mathbf{y}]E[\mathbf{y}]^T) &= \\ &= \sigma^2\text{Tr}(\mathbf{I}) - \sigma^2\text{Tr}(\mathbf{H}) + \text{Tr}((\mathbf{I} - \mathbf{H})E[\mathbf{y}]E[\mathbf{y}]^T) \quad (\because \text{property (2.5)}) \end{aligned}$$

Since $\mathbf{I} \in \mathbb{R}^{N \times 1}$, $\text{Tr}(\mathbf{I}) = N$. And $\text{Tr}(\mathbf{H}) = \text{Tr}(\mathbf{X}(\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T) = \text{Tr}((\mathbf{X}^T\mathbf{X})^{-1}(\mathbf{X}^T\mathbf{X}))$ because of property (2.8). So, $\text{Tr}(\mathbf{H}) = \text{Tr}(\mathbf{I}) = p + 1$ ($\because I \in \mathbb{R}^{(p+1) \times (p+1)}$)

$$\begin{aligned} \text{Tr}((\mathbf{I} - \mathbf{H})E[\mathbf{y}]E[\mathbf{y}]^T) &= \text{Tr}((\mathbf{I} - \mathbf{X}(\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T)(\mathbf{X}\beta)(\mathbf{X}\beta)^T) = \\ &= \text{Tr}((\mathbf{X}\beta - \mathbf{X}(\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T\mathbf{X}\beta)(\mathbf{X}\beta)^T) = \text{Tr}((\mathbf{X}\beta - \mathbf{X}\mathbf{I}\beta)(\mathbf{X}\beta)^T) = \\ &= \text{Tr}((\mathbf{X}\beta - \mathbf{X}\beta)(\mathbf{X}\beta)^T) = \text{Tr}(\mathbf{0}(\mathbf{X}\beta)^T), \quad \text{where } \mathbf{0} \text{ is a zero vector } \in \mathbb{R}^{N \times 1} \\ \text{Tr}(\mathbf{0}(\mathbf{X}\beta)^T) &= \text{Tr}(\mathbf{0}), \quad \mathbf{0} \text{ now represents a zero matrix } \in \mathbb{R}^{N \times N} \\ \text{Tr}(\mathbf{0}) &= 0 \quad (\text{a constant}) \end{aligned}$$

Therefore,

$$\begin{aligned} E[SSE] &= \sigma^2 \text{Tr}(\mathbf{I}) - \sigma^2 \text{Tr}(\mathbf{H}) + \text{Tr}((\mathbf{I} - \mathbf{H})E[\mathbf{y}]E[\mathbf{y}]^T) = \\ &= \sigma^2 N - \sigma^2(p+1) + 0 = \sigma^2(N-p-1) \end{aligned}$$

$$\text{So, } E[MSE] = E\left[\frac{SSE}{N-p-1}\right] = \frac{1}{N-p-1}E[SSE] = \frac{1}{N-p-1}\sigma^2(N-p-1) = \sigma^2$$

Exercise 3

For matrix \mathbf{A} , solve the following problems

$$\mathbf{A} = \begin{pmatrix} 6 & -3 \\ 5 & -2 \end{pmatrix}$$

3.1

Compute the eigenvalues λ_1, λ_2 ($\lambda_1 < \lambda_2$) and its corresponding eigenvectors v_1, v_2 of matrix \mathbf{A} [2 pts]

Answer:

Let's first compute eigenvalues.

$$\begin{aligned} \det(\mathbf{A} - \lambda \mathbf{I}) &= \det\left(\begin{matrix} 6-\lambda & -3 \\ 5 & -2-\lambda \end{matrix}\right) = (6-\lambda)(-2-\lambda) - (-3) \times 5 = \\ &= \lambda^2 - 4\lambda + 3 = 0 \Rightarrow \boxed{\lambda_1 = 3, \lambda_2 = 1} \end{aligned}$$

Now we will find eigenvectors.

$$\begin{aligned} \mathbf{A}v_1 &= \lambda_1 v_1 \quad \Rightarrow \quad \mathbf{A}v_1 - \lambda_1 v_1 = \bar{0}, \quad \bar{0} = \begin{pmatrix} 0 \\ 0 \end{pmatrix} \\ (\mathbf{A} - \lambda_1 \mathbf{I})v_1 &= \bar{0} \\ \begin{pmatrix} 6 - \lambda_1 & -3 \\ 5 & -2 - \lambda_1 \end{pmatrix}v_1 &= \begin{pmatrix} 6 - 3 & -3 \\ 5 & -2 - 3 \end{pmatrix}v_1 = \begin{pmatrix} 3 & -3 \\ 5 & -5 \end{pmatrix}v_1 = \bar{0} \end{aligned}$$

We will perform Gauss-Jordan elimination (but we will not follow the exact algorithm since the matrix is too simple) to find the vector $v_1 = \begin{pmatrix} v_{11} \\ v_{12} \end{pmatrix}$ and later the vector $v_2 = \begin{pmatrix} v_{21} \\ v_{22} \end{pmatrix}$:

$$\begin{pmatrix} 3 & -3 \\ 5 & -5 \end{pmatrix}v_1 = \bar{0} \Rightarrow \begin{pmatrix} 1 & -1 \\ 1 & -1 \end{pmatrix}v_1 = \bar{0} \Rightarrow \begin{pmatrix} 1 & -1 \\ 0 & 0 \end{pmatrix}v_1 = \bar{0} \Rightarrow v_{11} = v_{12}$$

We can freely choose a value for either v_{11} or v_{12} . Let's choose a value of 1 for v_{11} , which will mean that v_{12} will be equal to 1. So, $v_1 = \begin{pmatrix} 1 \\ 1 \end{pmatrix}$

Now, we will find v_2 .

$$\mathbf{A}v_2 = \lambda_2 v_2 \implies \mathbf{A}v_2 - \lambda_2 v_2 = \bar{0}, \quad \bar{0} = \begin{pmatrix} 0 \\ 0 \end{pmatrix}$$

$$(\mathbf{A} - \lambda_2 \mathbf{I})v_2 = \bar{0}$$

$$\begin{pmatrix} 6 - \lambda_2 & -3 \\ 5 & -2 - \lambda_2 \end{pmatrix} v_2 = \begin{pmatrix} 6 - 1 & -3 \\ 5 & -2 - 1 \end{pmatrix} v_2 = \begin{pmatrix} 5 & -3 \\ 5 & -3 \end{pmatrix} v_2 = \bar{0}$$

$$\begin{pmatrix} 5 & -3 \\ 5 & -3 \end{pmatrix} v_2 = \bar{0} \implies \begin{pmatrix} 1 & -0.6 \\ 1 & -0.6 \end{pmatrix} v_2 = \bar{0} \implies \begin{pmatrix} 1 & -0.6 \\ 0 & 0 \end{pmatrix} v_2 = \bar{0} \implies v_{21} = 0.6v_{22}$$

We can freely choose a value for either v_{21} or v_{22} . Let's choose a value of 1 for v_{22} , which will mean that v_{21} will be equal to 0.6. So, $v_2 = \begin{pmatrix} 0.6 \\ 1 \end{pmatrix}$

3.2

Find matrix \mathbf{P} to diagonalize \mathbf{A} . Here \mathbf{D} is a diagonal matrix of size 2×2 [3 pts]

$$\mathbf{P}^{-1}\mathbf{A}\mathbf{P} = \mathbf{D}$$

Answer:

$$\text{Let } \mathbf{V} = [v_1 \ v_2] = \begin{bmatrix} v_{11} & v_{21} \\ v_{12} & v_{22} \end{bmatrix}, \text{ and } \Lambda = \begin{bmatrix} \lambda_1 & 0 \\ 0 & \lambda_2 \end{bmatrix}.$$

Then the following equality holds:

$$\mathbf{AV} = \mathbf{V}\Lambda \tag{3.1}$$

Since the eigenvectors v_1 and v_2 that we found in problem 3.1 are independent, \mathbf{V}^{-1} exists, and thus we can multiply the both sides of the equation (3.1) by \mathbf{V}^{-1} from the left side:

$$\begin{aligned} \mathbf{V}^{-1}\mathbf{AV} &= \mathbf{V}^{-1}\mathbf{V}\Lambda \implies \\ &\implies \mathbf{V}^{-1}\mathbf{AV} = \mathbf{I}\Lambda = \Lambda \end{aligned} \tag{3.2}$$

Thus, matrix $\boxed{\mathbf{P} = \mathbf{V}}$

3.3

Compute the determinant of \mathbf{A}^{2521} . The calculation should be trivial if you use the properties of determinant. [3 pts]

Answer:

We are aware of the following property of the determinant when we have two square matrices \mathbf{A} and \mathbf{B} , both $\in \mathbb{R}^{n \times n}$:

$$\det(\mathbf{AB}) = \det(\mathbf{A})\det(\mathbf{B}) \tag{3.3}$$

The \mathbf{B} matrix can be \mathbf{A} also, so we can write the property (3.3) as

$$\begin{aligned} \det(\mathbf{A}^2) &= \det(\mathbf{A})\det(\mathbf{A}) \implies \\ \det(\mathbf{A}^3) &= \det(\mathbf{A}^2)\det(\mathbf{A}) = \det(\mathbf{A})\det(\mathbf{A})\det(\mathbf{A}) \end{aligned}$$

So, we can see the general trend as the power of \mathbf{A} grows, and thus we can write the following:

$$\det(\mathbf{A}^{2521}) = [\det(\mathbf{A})]^{2521}$$

Thus, if we can find the determinant of the matrix \mathbf{A} alone, we can find the determinant of \mathbf{A}^{2521} easily.

$$\det(\mathbf{A}) = \det\begin{pmatrix} 6 & -3 \\ 5 & -2 \end{pmatrix} = 6 \times (-2) - (-3) \times 5 = 3 \implies$$

$$\boxed{\det(\mathbf{A}^{2521}) = 3^{2521}}$$

Exercise 4

Given data \mathbf{X} , solve following problems

$$\mathbf{X} = \begin{pmatrix} 1 & 3 & 9 \\ 2 & 5 & 7 \\ 4 & 4 & 6 \\ 9 & 8 & 2 \end{pmatrix}$$

i.e., data with four samples and three features (predictors).

4.1

Find the mean value of each column. [1 pts]

Answer:

$$\begin{aligned}\mu_1 &= \frac{1 + 2 + 4 + 9}{4} = 4 \\ \mu_2 &= \frac{3 + 5 + 4 + 8}{4} = 5 \\ \mu_3 &= \frac{9 + 7 + 6 + 2}{4} = 6\end{aligned}$$

4.2

Subtract each mean from each element of the corresponding column.
Now, let us set the derived matrix as \mathbf{X}' . [1 pts]

Answer:

$$\mathbf{X}' = \begin{pmatrix} 1 - 4 & 3 - 5 & 9 - 6 \\ 2 - 4 & 5 - 5 & 7 - 6 \\ 4 - 4 & 4 - 5 & 6 - 6 \\ 9 - 4 & 8 - 5 & 2 - 6 \end{pmatrix} = \begin{pmatrix} -3 & -2 & 3 \\ -2 & 0 & 1 \\ 0 & -1 & 0 \\ 5 & 3 & -4 \end{pmatrix}$$

4.3

Calculate $\frac{1}{4-1} \mathbf{X}'^T \mathbf{X}'$. [1 pts]

Answer:

$$\begin{aligned}\frac{1}{4-1} \mathbf{X}'^T \mathbf{X}' &= \frac{1}{3} \begin{pmatrix} -3 & -2 & 3 \\ -2 & 0 & 1 \\ 0 & -1 & 0 \\ 5 & 3 & -4 \end{pmatrix}^T \begin{pmatrix} -3 & -2 & 3 \\ -2 & 0 & 1 \\ 0 & -1 & 0 \\ 5 & 3 & -4 \end{pmatrix} = \\ &= \frac{1}{3} \begin{pmatrix} -3 & -2 & 0 & 5 \\ -2 & 0 & -1 & 3 \\ 3 & 1 & 0 & -4 \end{pmatrix} \begin{pmatrix} -3 & -2 & 3 \\ -2 & 0 & 1 \\ 0 & -1 & 0 \\ 5 & 3 & -4 \end{pmatrix} = \frac{1}{3} \begin{pmatrix} 38 & 21 & -31 \\ 21 & 14 & -18 \\ -31 & -18 & 26 \end{pmatrix} = \begin{pmatrix} \frac{38}{3} & 7 & \frac{-31}{3} \\ 7 & \frac{14}{3} & -6 \\ \frac{-31}{3} & -6 & \frac{26}{3} \end{pmatrix}\end{aligned}$$

4.4

For the calculated matrix in 4.3, find eigenvalues $\lambda_1, \lambda_2, \lambda_3$ in descending order ($\lambda_1 > \lambda_2 > \lambda_3$) up to four decimal places. [2 pts]

(<https://www.symbolab.com/solver/matrix-eigenvalues-calculator/eigenvalues>)

Answer:

$$\lambda_1 = 25.3084$$

$$\lambda_2 = 0.6044$$

$$\lambda_3 = 0.0872$$

4.5

Calculate $\frac{\lambda_1}{(\lambda_1 + \lambda_2 + \lambda_3)}$. (up to four decimal places) [2 pts]

Answer:

$$\frac{\lambda_1}{(\lambda_1 + \lambda_2 + \lambda_3)} = \frac{25.3084}{25.3084 + 0.6044 + 0.0872} = 0.9734$$

4.6

Find eigenvectors $\mathbf{a}_1, \mathbf{a}_2, \mathbf{a}_3$ corresponding to $\lambda_1, \lambda_2, \lambda_3$. [2 pts]

(<https://www.symbolab.com/solver/matrix-eigenvalues-calculator/eigenvectors>)

Answer:

$$\mathbf{a}_1 = \begin{pmatrix} -1.2045 \\ -0.6992 \\ 1 \end{pmatrix}, \quad \mathbf{a}_2 = \begin{pmatrix} 138.5863 \\ -237.3327 \\ 1 \end{pmatrix}, \quad \mathbf{a}_3 = \begin{pmatrix} 0.6182 \\ 0.3652 \\ 1 \end{pmatrix}$$

Exercise 5

Consider three random variables X, Y, Z . The three variables have the covariance matrix in the form of:

$$\Sigma = \begin{pmatrix} a & ka & 0 \\ ka & a & ka \\ 0 & ka & a \end{pmatrix}$$

, where $0 < k < \frac{1}{\sqrt{2}}$

5.1

Calculate the eigenvalues $\lambda_1, \lambda_2, \lambda_3$. (Show process of calculation neatly) [4 pts]

Answer:

Date, No.

In order to find λ_1, λ_2 , and λ_3 , we should solve the follow. equation:

$$\det(\Sigma - \lambda I) = 0$$

$$\det(\Sigma - \lambda I) = \det \left(\begin{bmatrix} a & ka & 0 \\ ka & a & ka \\ 0 & ka & a \end{bmatrix} - \begin{bmatrix} \lambda & 0 & 0 \\ 0 & \lambda & 0 \\ 0 & 0 & \lambda \end{bmatrix} \right) =$$

$$= \det \left(\begin{bmatrix} a-\lambda & ka & 0 \\ ka & a-\lambda & ka \\ 0 & ka & a-\lambda \end{bmatrix} \right) =$$

$$= (a-\lambda) \det \left(\begin{bmatrix} a-\lambda & ka \\ ka & a-\lambda \end{bmatrix} \right) - ka \cdot \det \left(\begin{bmatrix} ka & ka \\ 0 & a-\lambda \end{bmatrix} \right)$$

$$+ 0 \cdot \det \left(\begin{bmatrix} ka & a-\lambda \\ 0 & ka \end{bmatrix} \right) = (a-\lambda)((a-\lambda)^2 - (ka)^2)$$

$$- ka((a-\lambda) - ka \cdot 0) + 0 = (a-\lambda)^3 - (ka)^2(a-\lambda)$$

$$- (ka)^2(a-\lambda) = 0 \Rightarrow (a-\lambda)^3 - 2(ka)^2(a-\lambda) = 0 \Rightarrow$$

$$(a-\lambda)((a-\lambda)^2 - 2(ka)^2) = 0 \Rightarrow$$

$$(a-\lambda)((a-\lambda) - \sqrt{2}(ka))((a-\lambda) + \sqrt{2}(ka)) = 0$$

$$\boxed{\lambda_1 = a}, \quad \boxed{\lambda_2 = a - \sqrt{2}ka}, \quad \boxed{\lambda_3 = a + \sqrt{2}ka}$$

5.2

Find PC(Principal Component) $\mathbf{p}_1, \mathbf{p}_2, \mathbf{p}_3$ of each random variable X, Y, Z . [3 pts]

Answer:

To find principal components, we should solve the follow.
equation: $(\Sigma - \lambda I) p = 0$ (zero vector)

For $\lambda_1 = \alpha$,

$$(\Sigma - \alpha I) p_1 = 0 \Rightarrow \begin{bmatrix} \alpha - Q & KQ & 0 & | & 0 \\ KQ & Q - Q & KQ & | & 0 \\ 0 & KQ & Q - Q & | & 0 \end{bmatrix} \rightarrow$$

$$p_1 = (p_{11} \ p_{12} \ p_{13})^T$$

$$\rightarrow \begin{bmatrix} 0 & KQ & 0 & | & 0 \\ KQ & 0 & KQ & | & 0 \\ 0 & KQ & 0 & | & 0 \end{bmatrix} \rightarrow p_{12} \times KQ = 0 \Rightarrow p_{11} KQ + p_{13} KQ = 0$$

$\left\{ \begin{array}{l} p_{12} \times KQ = 0 \\ (p_{11} + p_{13}) KQ = 0 \end{array} \right. \Rightarrow$ since $K > 0$ and Q represents the variance of a random variable, it cannot be zero. Thus we can divide the ~~two~~ equations by KQ , and then we get:

$$p_{12} = 0 \Rightarrow \text{If we choose } p_{13} = 1, \text{ then } p_1 = \begin{pmatrix} -1 \\ 0 \\ 1 \end{pmatrix}$$

$$p_{11} = -p_{13}$$

For $\lambda_2 = Q - \sqrt{2}KQ$,

$$(\Sigma - (Q - \sqrt{2}KQ) I) p_2 = 0 \Rightarrow \begin{bmatrix} Q - Q + \sqrt{2}KQ & KQ & 0 & | & 0 \\ KQ & Q - Q + \sqrt{2}KQ & KQ & | & 0 \\ 0 & KQ & Q - Q + \sqrt{2}KQ & | & 0 \end{bmatrix}$$

$$p_2 = (p_{21} \ p_{22} \ p_{23})^T$$

$$\rightarrow \begin{bmatrix} \sqrt{2}KQ & KQ & 0 & | & 0 \\ KQ & \sqrt{2}KQ & KQ & | & 0 \\ 0 & KQ & \sqrt{2}KQ & | & 0 \end{bmatrix} \rightarrow \begin{bmatrix} \sqrt{2} & 1 & 0 & | & 0 \\ 1 & \sqrt{2} & 1 & | & 0 \\ 0 & 1 & \sqrt{2} & | & 0 \end{bmatrix} \rightarrow$$

$$\rightarrow \left[\begin{array}{ccc|c} 1 & 1/\sqrt{2} & 0 & 0 \\ 1 & \sqrt{2} & 1 & 0 \\ 0 & 1 & \sqrt{2} & 0 \end{array} \right] \rightarrow \left[\begin{array}{ccc|c} 1 & 1/\sqrt{2} & 0 & 0 \\ 0 & 1/\sqrt{2} & 1 & 0 \\ 0 & 1 & \sqrt{2} & 0 \end{array} \right] \rightarrow$$

$$\left[\begin{array}{ccc|c} 1 & 1/\sqrt{2} & 0 & 0 \\ 0 & 1 & \sqrt{2} & 0 \\ 0 & 1 & \sqrt{2} & 0 \end{array} \right] \rightarrow \left[\begin{array}{ccc|c} 1 & 1/\sqrt{2} & 0 & 0 \\ 0 & 1 & \sqrt{2} & 0 \\ 0 & 0 & 0 & 0 \end{array} \right] \rightarrow$$

$$p_{22} = -\sqrt{2} p_{23}, \quad p_{21} = -\frac{p_{22}}{\sqrt{2}} = \frac{\sqrt{2} p_{23}}{\sqrt{2}} = 1$$

So, if we choose $p_{23} = 1$, then $p_2 = \begin{pmatrix} 1 \\ -\sqrt{2} \\ 1 \end{pmatrix}$

For $\lambda_3 = \theta + \sqrt{2} KQ$,

$$(\Sigma - (\theta + \sqrt{2} KQ) I) p_3 = 0 \Rightarrow \begin{bmatrix} \theta - \theta - \sqrt{2} KQ & KQ \\ KQ & \theta - \theta - \sqrt{2} KQ \\ 0 & KQ \end{bmatrix}$$

$$p_3 = (p_{31}, p_{32}, p_{33})^T$$

$$\left[\begin{array}{cc|c} 0 & 0 & 0 \\ KQ & 0 & 0 \\ 0 - \theta - \sqrt{2} KQ & 0 \end{array} \right] \Rightarrow \left[\begin{array}{ccc|c} -\sqrt{2} KQ & KQ & 0 & 0 \\ KQ & -\sqrt{2} KQ & KQ & 0 \\ 0 & KQ & -\sqrt{2} KQ & 0 \end{array} \right] \Rightarrow$$

$$\left[\begin{array}{ccc|c} -\sqrt{2} & 1 & 0 & 0 \\ 1 & -\sqrt{2} & 1 & 0 \\ 0 & 1 & -\sqrt{2} & 0 \end{array} \right] \Rightarrow \left[\begin{array}{ccc|c} 1 & -1/\sqrt{2} & 0 & 0 \\ 1 & -\sqrt{2} & 1 & 0 \\ 0 & 1 & -\sqrt{2} & 0 \end{array} \right] \Rightarrow$$

$$\left[\begin{array}{ccc|c} 1 & -1/\sqrt{2} & 0 & 0 \\ 0 & -1/\sqrt{2} & 1 & 0 \\ 0 & 1 & -\sqrt{2} & 0 \end{array} \right] \Rightarrow \left[\begin{array}{ccc|c} 1 & -1/\sqrt{2} & 0 & 0 \\ 0 & 1 & -\sqrt{2} & 0 \\ 0 & 1 & -\sqrt{2} & 0 \end{array} \right] \Rightarrow \left[\begin{array}{ccc|c} 1 & -1/\sqrt{2} & 0 & 0 \\ 0 & 1 & -\sqrt{2} & 0 \\ 0 & 0 & 0 & 0 \end{array} \right]$$

$p_{32} = \sqrt{2} p_{33}$, $p_{31} = \frac{1}{\sqrt{2}} p_{32} = p_{33}$. So, if we choose $p_{33} = 1$, then
 $p_3 = (1, \sqrt{2}, 1)^T$

5.3

Calculate how much total variance is explained by each principal component. [3 pts]

Answer:

$$\lambda_1 = Q, \lambda_2 = Q - \sqrt{2} KQ, \lambda_3 = Q + \sqrt{2} KQ$$

$$\text{Total variance} = \lambda_1 + \lambda_2 + \lambda_3 = 3Q$$

Let $V(p_i)$ represent the amount of variance explained by i -th principal component ($i = 1, 2, 3$)

$$V(p_1) = \frac{\lambda_1}{\lambda_1 + \lambda_2 + \lambda_3} = \frac{Q}{3Q} = \frac{1}{3}$$

$$V(p_2) = \frac{\lambda_2}{\lambda_1 + \lambda_2 + \lambda_3} = \frac{Q - \sqrt{2} KQ}{3Q} = \frac{1}{3} - \frac{\sqrt{2}}{3} K$$

$$V(p_3) = \frac{\lambda_3}{\lambda_1 + \lambda_2 + \lambda_3} = \frac{Q + \sqrt{2} KQ}{3Q} = \frac{1}{3} + \frac{\sqrt{2}}{3} K$$

So, we can see that p_3 explains highest

variance and after that p_1 and p_2 ($V(p_3) > V(p_1) > V(p_2)$)

Since, the problem did mention that we should

order eigenvalues and principal components, we

will leave things as they are.

Exercise 6

The following table is the outcome of the logistic regression model for an iris flower being species Versicolor versus species Virginica. ($Y = 1$ for Versicolor, $Y = 0$ for Virginica) [9 pts]

Variables	Intercept	Length of Sepals	Width of Sepals	Length of petals	Width of Petals
Estimated Coefficient	25.21	-3	-0.7	2.4	-10.3

6.1

Interpret the effect of length of sepals on the relative risks of an iris flower being species Versicolor versus species Virginica. Fill up the **A,B,C** in the below interpretation of the outcome and select the appropriate word for **D**. (Tips from TA : Think about the meaning of odds(X)!) [4 pts]

Answer:

Interpretation: Species **Virginica** is $\frac{1}{e^{-3}} = 20$ times more probable than Species **Versicolor** when the length of sepals **D increases** by 1 unit.

6.2

What is the predicted class for the following new data \mathbf{x}_1 and \mathbf{x}_2 ? You should provide the probability of $P(Y_1 = 1|\mathbf{x}_1), P(Y_1 = 0|\mathbf{x}_1), P(Y_2 = 1|\mathbf{x}_2), P(Y_2 = 0|\mathbf{x}_2)$. [6 pts]

\mathbf{x}_1	6.4	3.1	4.3	1.3
\mathbf{x}_2	6.9	3.0	3.9	1.4

Answer:

$$P(Y_1 = 1|\mathbf{x}_1) = \frac{1}{1 + e^{-(25.21+(-3)\times6.4+(-0.7)\times3.1+2.4\times4.3+(-10.3)\times1.3)}} = \frac{1}{1 + e^{-0.77}} = 0.6835$$

$$P(Y_1 = 0|\mathbf{x}_1) = 1 - P(Y_1 = 1|\mathbf{x}_1) = 1 - 0.6835 = 0.3165$$

So, the predicted class for \mathbf{x}_1 will be **Versicolor**.

$$P(Y_2 = 1|\mathbf{x}_2) = \frac{1}{1 + e^{-(25.21+(-3)\times6.9+(-0.7)\times3.0+2.4\times3.9+(-10.3)\times1.4)}} = \frac{1}{1 + e^{2.65}} = 0.066$$

$$P(Y_2 = 0|\mathbf{x}_2) = 1 - P(Y_2 = 1|\mathbf{x}_2) = 1 - 0.066 = 0.934$$

So, the predicted class for \mathbf{x}_2 will be **Virginica**.