

Ask for Help Print

OFM4 — OFM4 Task 2: Dimensionality Reduction Methods

Data Mining II — D212

PRFA — OFM4

TASK OVERVIEW

SUBMISSIONS

EVALUATION REPORT

Competencies

4030.6.5 : Dimension Reduction Methods

The graduate implements dimension reduction methods to identify significant variables.

Introduction

In this task, you will act as an analyst and create a data mining report. You must select one of the data dictionary and data set files to use for your report from the following web link: "[Data Sets and Associated Data Dictionaries](#)."

You should also refer to the data dictionary file for your chosen data set from the above link. You will use Python or R to analyze the given data and create a data mining report in a word processor (e.g., Microsoft Word). Throughout the submission, you must visually represent each step of your work and the findings of your data analysis.

Note: All algorithms and visual representations used need to be captured either in tables or as screenshots added into the submitted Word document. A separate Microsoft Excel (.xls or .xlsx) document of the cleaned data should be submitted along with the written aspects of the data mining report.

Scenario

Scenario 1

One of the most critical factors in customer relationship management that directly affects a company's long-term profitability is understanding the customers. When a company understands its customers' characteristics, it is better able to target products and marketing campaigns for them, resulting in better profits for the company in the long term.

You are an analyst for a telecommunications company that wants to better understand the characteristics of its customers. You have been asked to use principal component analysis (PCA) to analyze customer data to identify the principal variables of your customers, ultimately enabling better business and strategic decision-making.

Scenario 2

One of the most critical factors in patient relationship management that directly affects a hospital's long-term cost-effectiveness is understanding the patients and the conditions leading to hospital admissions. When a hospital understands its patients' characteristics, it is better able to target treatment to patients, resulting in a more effective cost of care for the hospital in the long term.

You are an analyst for a hospital that wants to better understand the characteristics of its patients. You have been asked to use PCA to analyze patient data to identify the principal variables of your patients, ultimately enabling better business and strategic decision-making for the hospital.

Requirements

Your submission must be your original work. No more than a combined total of 30% of the submission and no more than a 10% match to any one individual source can be directly quoted or closely paraphrased from sources, even if cited correctly. The similarity report that is provided when you submit your task can be used as a guide.

You must use the rubric to direct the creation of your submission because it provides detailed criteria that will be used to evaluate your work. Each requirement below may be evaluated by more than one rubric aspect. The rubric aspect titles may contain hyperlinks to relevant portions of the course.

*Tasks may **not** be submitted as cloud links, such as links to Google Docs, Google Slides, OneDrive, etc., unless specified in the task requirements. All other submissions must be file types that are uploaded and submitted as attachments (e.g., .docx, .pdf, .ppt).*

Part I: Research Question

A. Describe the purpose of your data mining report by doing the following:

1. Propose **one** question relevant to a real-world organizational situation that you will answer by using PCA.
2. Define **one** goal of the data analysis. Ensure your goal is reasonable within the scope of the selected scenario and is represented in the available data.

Part II: Method Justification

B. Explain the reasons for using PCA by doing the following:

1. Explain how PCA analyzes the selected data set. Include expected outcomes.
2. Summarize **one** assumption of PCA.

Part III: Data Preparation

C. Perform data preparation for the chosen data set by doing the following:

1. Identify the continuous data set variables that you will need to answer the PCA question proposed in part A1.
2. Standardize the continuous data set variables identified in part C1. Include a copy of the cleaned data set.

Part IV: Analysis

D. Perform PCA by doing the following:

1. Determine the matrix of *all* the principal components.
2. Identify the *total* number of principal components, using the elbow rule or the Kaiser criterion. Include a screenshot of the scree plot.
3. Identify the variance of *each* of the principal components identified in part D2.
4. Identify the *total* variance captured by the principal components identified in part D2.

5. Summarize the results of your data analysis.

Part V: Attachments

E. Record the web sources used to acquire data or segments of third-party code to support the analysis. Ensure the web sources are reliable.

F. Acknowledge sources, using in-text citations and references, for content that is quoted, paraphrased, or summarized.

G. Demonstrate professional communication in the content and presentation of your submission.

File Restrictions

File name may contain only letters, numbers, spaces, and these symbols: ! - _ . * ' ()

File size limit: 200 MB

File types allowed: doc, docx, rtf, xls, xlsx, ppt, pptx, odt, pdf, csv, txt, qt, mov, mpg, avi, mp3, wav, mp4, wma, flv, asf, mpeg, wmv, m4v, svg, tif, tiff, jpeg, jpg, gif, png, zip, rar, tar, 7z

Rubric

A1:PROPOSAL OF QUESTION

NOT EVIDENT

The submission does not propose 1 question.

APPROACHING COMPETENCE

The submission proposes 1 question that cannot be answered using PCA, or the question is not relevant to a real-world organizational situation.

COMPETENT

The submission proposes 1 question that can be answered using PCA and is relevant to a real-world organizational situation.

A2:DEFINED GOAL

NOT EVIDENT

The submission does not define 1 goal of the data analysis.

APPROACHING COMPETENCE

The submission defines 1 goal of the data analysis, but the goal is not reasonable within the scope of the selected scenario or is not represented in the available data.

COMPETENT

The submission defines 1 reasonable goal of the data analysis that is within the scope of the scenario and is represented in the available data.

B1:EXPLANATION OF PCA

NOT EVIDENT

APPROACHING COMPETENCE

COMPETENT

The submission does not explain how PCA analyzes the selected data set.

The submission does not logically explain how PCA analyzes the selected data set or does not include expected outcomes.

The submission logically explains how PCA analyzes the selected data set and includes expected outcomes.

B2:PCA ASSUMPTION

NOT EVIDENT

The submission does not summarize 1 assumption of PCA.

APPROACHING COMPETENCE

The submission inaccurately summarizes 1 assumption of PCA.

COMPETENT

The submission accurately summarizes 1 assumption of PCA.

C1:CONTINUOUS DATA SET VARIABLES

NOT EVIDENT

The submission does not identify the continuous data set variables needed to answer the PCA question from part A1.

APPROACHING COMPETENCE

The submission inaccurately identifies the continuous data set variables needed to answer the PCA question from part A1.

COMPETENT

The submission accurately identifies the continuous data set variables needed to answer the PCA question from part A1.

C2:STANDARDIZATION OF DATA SET VARIABLES

NOT EVIDENT

The submission does not standardize the continuous data set variables.

APPROACHING COMPETENCE

The submission inaccurately standardizes the continuous data set variables identified in part C1 or does not include a cleaned data set.

COMPETENT

The submission accurately standardizes the continuous data set variables identified in part C1 and includes a cleaned data set.

D1:PRINCIPAL COMPONENTS

NOT EVIDENT

The submission does not determine the matrix of the principal components.

APPROACHING COMPETENCE

The submission inaccurately determines the matrix of 1 or more of the principal components.

COMPETENT

The submission accurately determines the matrix of *all* the principal components.

D2:IDENTIFICATION OF THE TOTAL NUMBER OF COMPONENTS**NOT EVIDENT**

The submission does not identify the *total* number of principal components.

APPROACHING COMPETENCE

The submission inaccurately identifies the *total* number of principal components, or it does not use the elbow rule or the Kaiser criterion or does not include a screenshot of the scree plot.

COMPETENT

The submission accurately identifies the *total* number of principal components using the elbow rule or the Kaiser criterion. The submission includes a screenshot of the scree plot.

D3:VARIANCE OF EACH COMPONENT**NOT EVIDENT**

The submission does not identify the variance of *any* of the principal components identified in part D2.

APPROACHING COMPETENCE

The submission inaccurately identifies the variance of 1 or more of the principal components identified in part D2.

COMPETENT

The submission accurately identifies the variance of *each* of the principal components identified in part D2.

D4:TOTAL VARIANCE CAPTURED BY COMPONENTS**NOT EVIDENT**

The submission does not identify the *total* variance captured by the principal components identified in part D2.

APPROACHING COMPETENCE

The submission inaccurately identifies the *total* variance captured by the principal components identified in part D2.

COMPETENT

The submission accurately identifies the *total* variance captured by the principal components identified in part D2.

D5:SUMMARY OF DATA ANALYSIS

NOT EVIDENT

The submission does not summarize the results of the data analysis.

APPROACHING COMPETENCE

The submission inaccurately summarizes the results of the data analysis.

COMPETENT

The submission accurately summarizes the results of the data analysis.

E:SOURCES FOR THIRD-PARTY CODE

NOT EVIDENT

The submission does not record web sources used to acquire data or segments of third-party code.

APPROACHING COMPETENCE

The submission records 1 or more unreliable web sources used to acquire data or segments of third-party code.

COMPETENT

The submission records *all* web sources used to acquire data or segments of third-party code, and the web sources are reliable.

F:SOURCES

NOT EVIDENT

The submission does not include both in-text citations and a reference list for sources that are quoted, paraphrased, or summarized.

APPROACHING COMPETENCE

The submission includes in-text citations for sources that are quoted, paraphrased, or summarized and a reference list; however, the citations or reference list is incomplete or inaccurate.

COMPETENT

The submission includes in-text citations for sources that are properly quoted, paraphrased, or summarized and a reference list that accurately identifies the author, date, title, and source location as available.

G:PROFESSIONAL COMMUNICATION

NOT EVIDENT

Content is unstructured, is disjointed, or contains pervasive

APPROACHING COMPETENCE

Content is poorly organized, is difficult to follow, or contains errors

COMPETENT

Content reflects attention to detail, is organized, and focuses on the main ideas as prescribed in the task or

errors in mechanics, usage, or grammar. Vocabulary or tone is unprofessional or distracts from the topic.

in mechanics, usage, or grammar that cause confusion. Terminology is misused or ineffective.

chosen by the candidate. Terminology is pertinent, is used correctly, and effectively conveys the intended meaning. Mechanics, usage, and grammar promote accurate interpretation and understanding.

Web Links

Data Sets and Associated Data Dictionaries

If you have trouble with the link, copy and paste the URL directly into your web browser.