

Data Mining II
PA2 Dimensionality Reduction Methods

Jason Willis

College of Information Technology,

Western Governors University

Dr. Kesselly Kamara

June 25th, 2022

Table of Contents for Each Rubric***Part I: Research Question******Describe Purpose, Summarize Research Question and Define Objectives:*** 3***Part II: Method Justification******Explain PCA Analysis & Summarize Assumptions*** 4***Part III: Data Preparation******Perform Data Preparation, Identify and Standardize Continuous data*** 6***Part IV: Perform PCA Analysis******Determine matrix, Identify All Principal Components and Variance, Summarize*** 8***Part V: Attachments******Sources for Third-Party Code*** 10***Sources*** 11

Hospital Readmission Problem

For our chain of hospitals to lower readmission concerns, we need to identify patients who have increased risk of rehospitalization within a month of their release. According to Schuller (2020), non-obese adults were 21% less likely to be readmitted than obese adults. A readmission study by Gert, et. al. (2002) showed a correlation between longer initial hospital stays and readmission. Within the provided dataset, I'm leveraging these studies to help create my hypothetical question and shape my approach in finding potential patient groups with a statistically significant chance for readmission outcomes.

After viewing the provided medical_clean.csv data set and accompanying data dictionary, there seems to be some patient groupings which are aligned with the research mentioned above. For instance, the following patient data fields: Initial patient admin days, Total Charges, and Initial Says (inpatient) both caught my attention and were underscored by the research mentioned above. While my initial feelings towards these variables might make them feel related, are they?

A1 – Proposal of Question

From information about previous patients who were readmitted, can we ascertain the minimum number of principal variables for our patients?

A2 – Defined Goal

The goal of our analysis is to logically investigate the provided data set and, by leveraging principal component analysis and other techniques, reduce the dimensionality of the provided dataset while keeping the model's accuracy high.

B1 – Explanation of PCA

According to Larose (2019) “PCA seeks to account for the correlation structure of a set of predictor variables, using a smaller set of uncorrelated linear combinations of these variables, called components.” The approach aims to reduce the original number of predictor variables while accounting for most predictions in the original set. By reducing the predictor set, redundancies like multicollinearity are removed. Additionally, reductions in predictor values also helps to decrease overfitting.

B2 – PCA Assumption:

As stated in the Principal Component Analysis (PCA) explanation above, the technique assumes reduced data set k can provide almost indistinguishable predictive results as compared to the original complete data set m .

C1 – Continuous Dataset Variables

The data set accurately identifies continuous (and discrete) dataset variables needed to answer the PCA question from part A1 (Figure 1, Figure 2).

| Create DataFrame w/Number DataTypes Only | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
|--|--------------------------------------|----------------------------------|----------|------------|----------|-----|--------|----------|--|-------------|-------------|-----|------|------------|-----------|--------------------------------------|----------------------------------|-----|--|-----------|--------------------------------------|----------------------------------|----------|-----------|------|---|----|----------|--|---|---|-------|----------|-----------|-------|---|----|----------|--|---|---|-------|----------|-----------|-------|---|----|----------|--|---|---|-------|----------|-----------|------|---|----|----------|--|---|---|-------|----------|-----------|------|---|----|---------|--|
| <pre>1 df_num = df.select_dtypes(include='number') 2 df_num.head()</pre> | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| <table border="1"> <thead> <tr> <th>CaseOrder</th><th>Zip</th><th>Lat</th><th>Lng</th><th>Population</th><th>Children</th><th>Age</th><th>Income</th><th></th><th></th></tr> </thead> <tbody> <tr> <td>0</td><td>1</td><td>35621</td><td>34.34960</td><td>-86.72508</td><td>2951</td><td>1</td><td>53</td><td>86575.93</td><td></td></tr> <tr> <td>1</td><td>2</td><td>32446</td><td>30.84513</td><td>-85.22907</td><td>11303</td><td>3</td><td>51</td><td>46805.99</td><td></td></tr> <tr> <td>2</td><td>3</td><td>57110</td><td>43.54321</td><td>-96.63772</td><td>17125</td><td>3</td><td>53</td><td>14370.14</td><td></td></tr> <tr> <td>3</td><td>4</td><td>56072</td><td>43.89744</td><td>-93.51479</td><td>2162</td><td>0</td><td>78</td><td>39741.49</td><td></td></tr> <tr> <td>4</td><td>5</td><td>23181</td><td>37.59894</td><td>-76.88958</td><td>5287</td><td>1</td><td>22</td><td>1209.56</td><td></td></tr> </tbody> </table> | | | | | | | | | | CaseOrder | Zip | Lat | Lng | Population | Children | Age | Income | | | 0 | 1 | 35621 | 34.34960 | -86.72508 | 2951 | 1 | 53 | 86575.93 | | 1 | 2 | 32446 | 30.84513 | -85.22907 | 11303 | 3 | 51 | 46805.99 | | 2 | 3 | 57110 | 43.54321 | -96.63772 | 17125 | 3 | 53 | 14370.14 | | 3 | 4 | 56072 | 43.89744 | -93.51479 | 2162 | 0 | 78 | 39741.49 | | 4 | 5 | 23181 | 37.59894 | -76.88958 | 5287 | 1 | 22 | 1209.56 | |
| CaseOrder | Zip | Lat | Lng | Population | Children | Age | Income | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| 0 | 1 | 35621 | 34.34960 | -86.72508 | 2951 | 1 | 53 | 86575.93 | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| 1 | 2 | 32446 | 30.84513 | -85.22907 | 11303 | 3 | 51 | 46805.99 | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| 2 | 3 | 57110 | 43.54321 | -96.63772 | 17125 | 3 | 53 | 14370.14 | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| 3 | 4 | 56072 | 43.89744 | -93.51479 | 2162 | 0 | 78 | 39741.49 | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| 4 | 5 | 23181 | 37.59894 | -76.88958 | 5287 | 1 | 22 | 1209.56 | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| 5 rows × 23 columns | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| ••• | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| ••• | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| Create DataFrame w/Categorical DataTypes Only | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| <pre>1 df_cat = df.select_dtypes(exclude='number') 2 df_cat.head()</pre> | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| <table border="1"> <thead> <tr> <th>Customer_id</th><th>Interaction</th><th>UID</th><th>City</th><th>St</th></tr> </thead> <tbody> <tr> <td>0 C412403</td><td>8cd49b13-f45a-4b47-a2bd-173ffa932c2f</td><td>3a83ddb66e2ae73798bdf1d705dc0932</td><td>Eva</td><td></td></tr> <tr> <td>1 Z919181</td><td>d2450b70-0337-4406-bdbb-bc1037f1734c</td><td>176354c5eef714957d486009feabf195</td><td>Marianna</td><td></td></tr> </tbody> </table> | | | | | | | | | | Customer_id | Interaction | UID | City | St | 0 C412403 | 8cd49b13-f45a-4b47-a2bd-173ffa932c2f | 3a83ddb66e2ae73798bdf1d705dc0932 | Eva | | 1 Z919181 | d2450b70-0337-4406-bdbb-bc1037f1734c | 176354c5eef714957d486009feabf195 | Marianna | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| Customer_id | Interaction | UID | City | St | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| 0 C412403 | 8cd49b13-f45a-4b47-a2bd-173ffa932c2f | 3a83ddb66e2ae73798bdf1d705dc0932 | Eva | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| 1 Z919181 | d2450b70-0337-4406-bdbb-bc1037f1734c | 176354c5eef714957d486009feabf195 | Marianna | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |

Figure 1 - Identifying Continuous and Discrete Variables

```
[159]: 1 pca_df_binary_standardized.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 10000 entries, 0 to 9999
Data columns (total 40 columns):
 #   Column           Non-Null Count  Dtype  
--- 
 0   Lat              10000 non-null   float64
 1   Lng              10000 non-null   float64
 2   Population       10000 non-null   float64
 3   Children         10000 non-null   float64
 4   Age              10000 non-null   float64
 5   Income            10000 non-null   float64
 6   VitD_levels      10000 non-null   float64
 7   Doc_visits       10000 non-null   float64
 8   Full_meals_eaten 10000 non-null   float64
 9   vitD_supp        10000 non-null   float64
 10  Initial_days     10000 non-null   float64
 11  TotalCharge      10000 non-null   float64
 12  Additional_charges 10000 non-null   float64
 13  Area_Suburban    10000 non-null   float64
 14  Area_Urban       10000 non-null   float64
 15  Marital_Married  10000 non-null   float64
 16  Marital_Never Married 10000 non-null   float64
 17  Marital_Separated 10000 non-null   float64
 18  Marital_Widowed  10000 non-null   float64
 19  Gender_Male      10000 non-null   float64
 20  Gender_Nonbinary 10000 non-null   float64
 21  Soft_drink_Yes   10000 non-null   float64
 22  Initial_admin_Emergency Admission 10000 non-null   float64
 23  Initial_admin_Observation Admission 10000 non-null   float64
 24  HighBlood_Yes    10000 non-null   float64
 25  Stroke_Yes       10000 non-null   float64
 26  Complication_risk_Low 10000 non-null   float64
 27  Complication_risk_Medium 10000 non-null   float64
 28  Overweight_Yes   10000 non-null   float64
 29  Arthritis_Yes    10000 non-null   float64
 30  Diabetes_Yes     10000 non-null   float64
 31  Hyperlipidemia_Yes 10000 non-null   float64
 32  BackPain_Yes     10000 non-null   float64
 33  Anxiety_Yes      10000 non-null   float64
 34  Allergic_rhinitis_Yes 10000 non-null   float64
 35  Reflux_esophagitis_Yes 10000 non-null   float64
 36  Asthma_Yes       10000 non-null   float64
 37  Services_CT Scan 10000 non-null   float64
 38  Services_Intravenous 10000 non-null   float64
 39  Services_MRI      10000 non-null   float64
dtypes: float64(40)
memory usage: 3.1 MB
```

Figure 2 - Data Types

C2 – Standardization of Dataset Variables

During the preprocessing steps, unnecessary features were removed dummy variables (Figure 3), were created.

```
[52]: 1 df_temp = df[['Age', 'Gender', 'ReAdmis', 'VitD_levels', 'Doc_visits', 'vitD_supp', 'Initial_admin', \
2   'HighBlood', 'Stroke', 'Complication_risk', 'Overweight', 'Arthritis', 'Diabetes', 'Hyperlipidemia', \
3   'BackPain', 'Anxiety', 'Allergic_rhinitis', 'Reflux_esophagitis', 'Asthma', 'Services', 'Initial_days', \
4   'TotalCharge', 'Additional_charges']]
```

```
[56]: 1 df_dummies = pd.get_dummies(df_temp, drop_first=True)
2 df_dummies.head()
```

| | Age | VitD_levels | Doc_visits | vitD_supp | Initial_days | TotalCharge | Additional_charges | Gender_Male | Gender_Nonbinary | ReAdmis_Yes | ... | Diabetes_Yes | Hyperlipidemia_Yes | ... |
|---|-----|-------------|------------|-----------|--------------|-------------|--------------------|-------------|------------------|-------------|-----|--------------|--------------------|-----|
| 0 | 53 | 19.141466 | 6 | 0 | 10.585770 | 3726.702860 | 17939.403420 | 1 | 0 | 0 | ... | 1 | 0 | 0 |
| 1 | 51 | 18.940352 | 4 | 1 | 15.129662 | 4193.190458 | 17612.998120 | 0 | 0 | 0 | ... | 0 | 0 | 0 |
| 2 | 53 | 18.057507 | 4 | 0 | 4.772177 | 2434.234222 | 17505.192460 | 0 | 0 | 0 | ... | 1 | 0 | 0 |
| 3 | 78 | 16.578658 | 4 | 0 | 1.744879 | 2127.830423 | 12993.437350 | 1 | 0 | 0 | ... | 0 | 0 | 0 |
| 4 | 22 | 17.439069 | 5 | 2 | 1.254807 | 2113.073274 | 3716.525786 | 0 | 0 | 0 | ... | 0 | 0 | 1 |

5 rows × 28 columns

Figure 3 - Create Dummy Data

Additionally, the dataset was standardized using `sklearn.preprocessing StandardScaler()` function (Figure 4) and then exported as `std_df.csv`.

Calculating Principle Components

```
1 from sklearn.preprocessing import StandardScaler  
2 from sklearn.decomposition import PCA  
3  
4 # Standardizing Continuous Dataset  
5 scaler = StandardScaler()  
6 std_df = scaler.fit_transform(df_dummies)  
7  
8 pca = PCA()  
9 print(pca.fit_transform(std_df))  
  
[ [-1.22345671e+00  1.38470718e+00 -8.99594122e-01 ... -8.02858050e-02  
  2.95935943e-02 -4.74054638e-01]  
[-1.07180686e+00  1.05141455e+00 -1.10734035e+00 ...  5.92369348e-02  
  5.00476368e-02 -1.58767282e-09]  
[-1.94004422e+00  1.10048864e+00  5.73019886e-01 ... -2.23728403e-01  
  2.00594989e-02  3.79674304e-09]  
...  
[ 2.26175452e+00  4.45697534e-01  4.88310313e-01 ... -8.21186227e-02  
  9.77371257e-02  2.67002221e-08]  
[ 1.97986293e+00 -1.19281366e+00 -7.54375873e-01 ... -1.68552210e-01  
 -9.31021652e-02  1.12305098e-07]  
[ 2.16046478e+00 -4.67578843e-01  1.54701307e+00 ... -5.68823370e-02  
  7.47586277e-02  2.00799114e-07]]
```

Export Cleaned and Standardized Dataset

```
1 pd.DataFrame(std_df).to_csv('std_df.csv', index=False)
```

Figure 4 - Standardize and Export Continuous Dataset

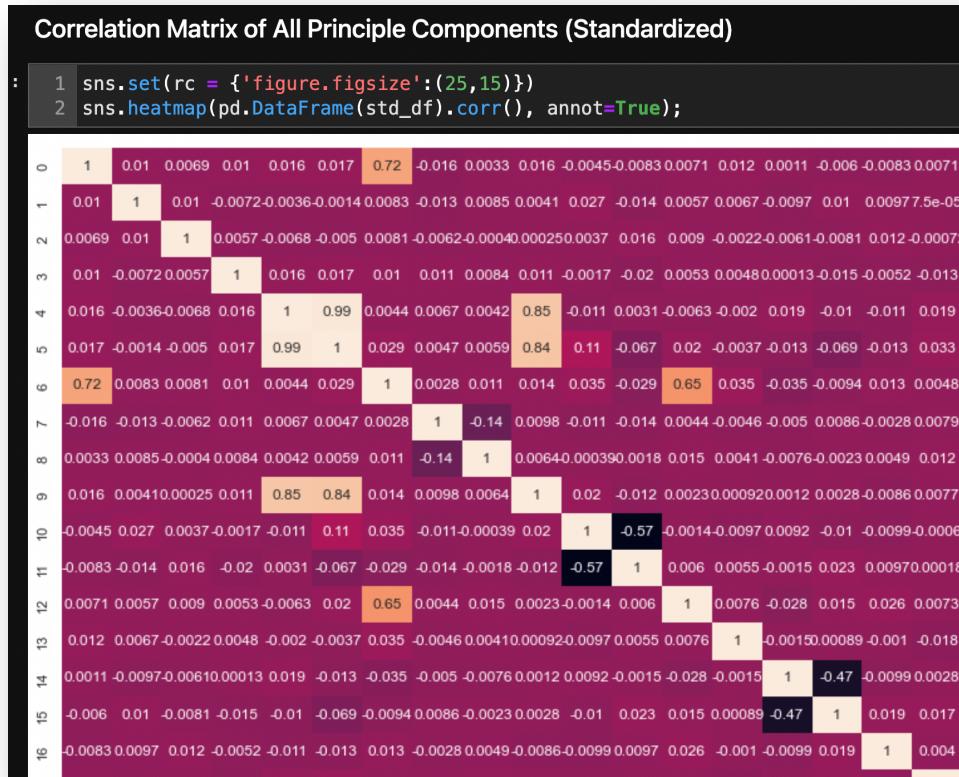


Figure 5 - Matrix of All Principle Components

D1 – Principal Components

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

(Figure 5).

D2 – Identification of Total Number of Components

The submission accurately identifies the total number of principal components. (Figure 6). The ideal number of components are 3 with an accuracy of 95.5%.

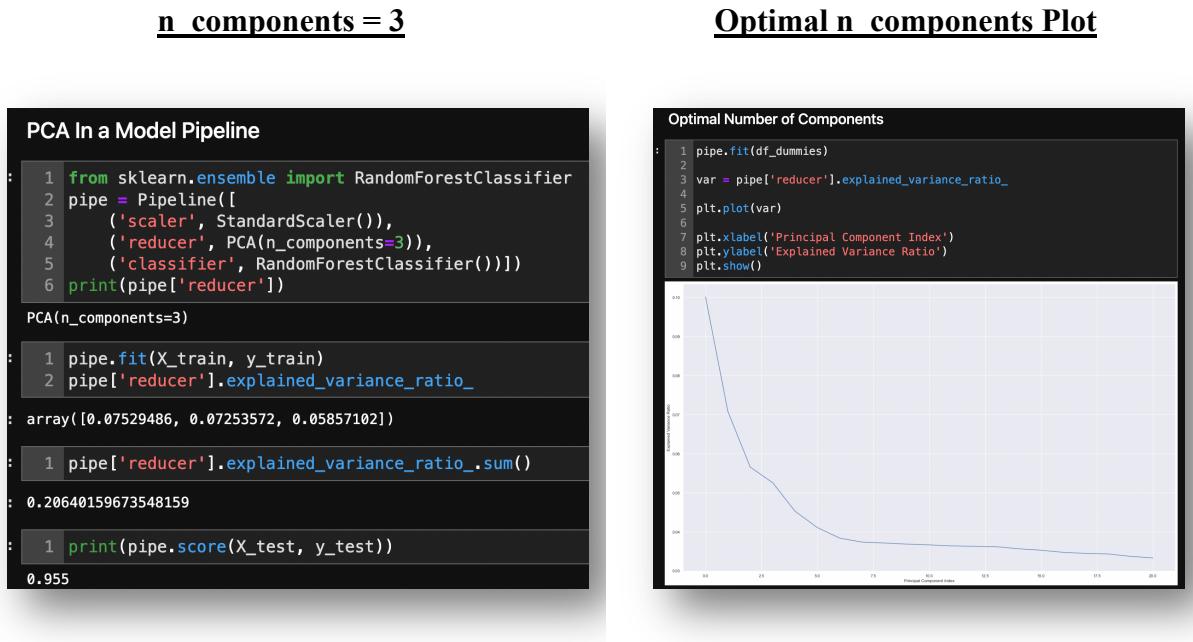


Figure 6 - Optimal n_components

D3 – Total Variance of Components

The submission accurately identifies the variance of each of the principal components identified in part D2 (Figure 6). You can see that each principal component's variance ratio and their sum.

D4 – Total Variance Captured by Components

The submission accurately identifies the total variance captured by the principal components identified in part D2 using the explained_variance_ratio_.sum(). (Figure 6)

D5 – Summary of Data Analysis

At 95.5% accuracy, this model is quite precise while reducing principal component features. I would have a high confidence in the dimension reduction process taken here to prepare the dataset. I believe we achieved our goal of figuring out the minimum number of principal variables for our patients while maintaining the model's accuracy high.

E – Sources for Third-Party Code

- Help using Markdown: <https://www.markdownguide.org/basic-syntax/>
- Help to see ALL columns: <https://stackoverflow.com/questions/24524104/pandas-describe-is-not-returning-summary-of-all-columns>
- Help to create a better histogram design: https://mode.com/example-gallery/python_histogram/
- Matplotlib Help: https://matplotlib.org/2.1.2/api/_as_gen/matplotlib.pyplot.plot.html
- Multiple ways to conduct ANOVA: <https://www.marsja.se/four-ways-to-conduct-one-way-anovas-using-python/>
- Numpy Help: <https://numpy.org/doc/stable/>
- Pandas Help: https://pandas.pydata.org/docs/user_guide/index.html#user-guide
- Python Help: <https://docs.python.org/3.9/library/index.html>
- Scipy.stats Help: <https://docs.scipy.org/doc/scipy/reference/tutorial/stats.html>

References

- Gert P Westert, Ronald J Lagoe, Ilmo Keskimäki, Alastair Leyland, Mark Murphy,
An international study of hospital readmissions and related utilization in Europe and the
USA, Health Policy, Volume 61, Issue 3, 2002, Pages 269-278, ISSN 0168-8510,
[https://doi.org/10.1016/S0168-8510\(01\)00236-6](https://doi.org/10.1016/S0168-8510(01)00236-6).
(<https://www.sciencedirect.com/science/article/pii/S0168851001002366>)
- Larose, D., C., & Larose, D., T. (2019). Data Science Using Python and R. Wiley.
<https://www.wiley.com/en-us/Data+Science+Using+Python+and+R-p-9781119526810>
- Schuller K. A. (2020). Is obesity a risk factor for readmission after acute myocardial
infarction? *Journal of healthcare quality research*, 35(1), 4–11.
<https://doi.org/10.1016/j.jhqr.2019.09.002>