

# Feature Extraction

Ruwen Qin  
Stony Brook University

September 4, 2025

Data used in this module:

- California house prices dataset from `sklearn.datasets`
- Optical recognition of handwritten digits dataset from `sklearn.datasets`

Python notebook used in this chapter

- `Feature_Extraction.ipynb`

# 1 Background

Considering a system characterized by  $N$  attributes. We collect a sample data of the attributes, which contains  $M$  observations. Denote the sample data as  $\mathbf{X} \in \mathbb{R}^{M \times N}$ . If  $N$  is large, we hope to characterize the system with fewer but informative features. In some cases  $N$  is not large, but some attributes may be correlated with each other. We hope to find features that independently characterize the system from unique perspectives.

In this module, we will study principal component analysis and then briefly introduce a few more ML methods for feature extraction and dimension reduction.

## 2 Principal Component Analysis (PCA)

Define  $\mathcal{P}_d$  as the set of  $N$ -dimensional rank- $d$  orthogonal projection matrix<sup>1</sup>. We can project the zero-centered sample data  $\tilde{\mathbf{X}} (= \mathbf{X} - \bar{\mathbf{X}})$  onto the  $d$ -dimensional linear sub-space using a projection matrix  $\mathbf{P} \in \mathcal{P}_d$ . We hope to minimize the reconstruction error; that is, the sum of  $L_2$  distances between the original data  $\tilde{\mathbf{X}}$  and the projected data  $\tilde{\mathbf{X}}\mathbf{P}$ :

$$\min_{\mathbf{P} \in \mathcal{P}_d} \|\tilde{\mathbf{X}}\mathbf{P} - \tilde{\mathbf{X}}\|_F^2, \quad (1)$$

where  $\|\cdot\|_F$  is the Frobenius norm for matrices.

The *Principal Component Analysis* (PCA) [1] shows that the solution to the minimization problem in (1),  $\mathbf{P}^*$ , contains  $d$  columns that are top  $d$  eigenvectors of  $\mathbf{X}$ 's covariance matrix  $\Sigma$ :

$$\Sigma = \frac{(\mathbf{X} - \bar{\mathbf{X}})^T(\mathbf{X} - \bar{\mathbf{X}})}{N - 1} = \frac{\tilde{\mathbf{X}}^T \tilde{\mathbf{X}}}{N - 1} \quad (2)$$

with eigenvalues  $\lambda_1 \geq \dots \geq \lambda_d$ <sup>2</sup>.

### 2.1 Principal Components

Principal components (PCs) are normalized linear combinations of the feature vectors:

$$\mathbf{C}_j = \tilde{\mathbf{X}}\mathbf{p}_j \quad (3)$$

The variance of the  $j$ th PC is

$$\text{Var}[\mathbf{C}_j] = \mathbf{p}_j^T \Sigma \mathbf{p}_j, \quad (4)$$

which is the eigenvalue  $\lambda_j$  corresponding to the eigenvector  $\mathbf{p}_j$ . Therefore,

$$\frac{\lambda_j}{\sum_{i=1}^d \lambda_i} \quad (5)$$

is the proportion of variance explained by the  $j$ th PC and

$$\frac{\sum_{i=1}^j \lambda_i}{\sum_{i=1}^d \lambda_i} \quad (6)$$

is the proportion of variance explained by the first  $j$  PCs.

---

<sup>1</sup>an  $N \times N$  square matrix that can be written as  $UU^T$  where  $U \in \mathbb{R}^{N \times d}$  has  $d$  orthogonal columns.

<sup>2</sup>In linear algebra, we also know that the eigenvectors of  $\mathbf{X}$ 's co-variance matrix are the right singular vectors of  $\tilde{\mathbf{X}}$  and the eigenvalues of  $\mathbf{X}$ 's co-variance matrix are squares of  $\tilde{\mathbf{X}}$ 's positive eigenvalues divided by  $(N-1)$

## 2.2 Derivation of PCA Algorithm

To minimize the reconstruction loss

$$\begin{aligned}\|\tilde{\mathbf{X}}\mathbf{P} - \tilde{\mathbf{X}}\|_F^2 &= \text{Tr}((\tilde{\mathbf{X}}\mathbf{P} - \tilde{\mathbf{X}})(\tilde{\mathbf{X}}\mathbf{P} - \tilde{\mathbf{X}})^T) \\ &= \text{Tr}(\tilde{\mathbf{X}}\mathbf{P}\mathbf{P}^T\tilde{\mathbf{X}}^T - 2\tilde{\mathbf{X}}\mathbf{P}\tilde{\mathbf{X}}^T + \tilde{\mathbf{X}}\tilde{\mathbf{X}}^T) \\ &= -\text{Tr}(\tilde{\mathbf{X}}\mathbf{P}\tilde{\mathbf{X}}^T) + \text{Tr}(\tilde{\mathbf{X}}\tilde{\mathbf{X}}^T)\end{aligned}\quad (7)$$

is to maximize  $\text{Tr}(\tilde{\mathbf{X}}\mathbf{P}\tilde{\mathbf{X}}^T)$ <sup>3</sup>.

Because  $\mathbf{P}$  is an  $N$ -dimensional rank- $d$  orthogonal projection matrix,  $\mathbf{P} = \mathbf{U}\mathbf{U}^T$ , where  $\mathbf{U} \in \mathbb{R}^{N \times d}$  containing orthogonal columns. Therefore,

$$\text{Tr}(\tilde{\mathbf{X}}\mathbf{P}\tilde{\mathbf{X}}^T) = \text{Tr}(\tilde{\mathbf{X}}\mathbf{U}\mathbf{U}^T\tilde{\mathbf{X}}^T) = \text{Tr}(\mathbf{U}^T\tilde{\mathbf{X}}^T\tilde{\mathbf{X}}\mathbf{U})^4 = \sum_{j=1}^d \mathbf{u}_j^T \tilde{\mathbf{X}}^T \tilde{\mathbf{X}} \mathbf{u}_j, \quad (8)$$

where  $\mathbf{u}_j$  is the  $j$ th column of  $\mathbf{U}$ .

Therefore, we solve the following optimization problem:

$$\begin{aligned}\max_{\mathbf{u}_1, \dots, \mathbf{u}_d} \quad & \sum_{j=1}^d \mathbf{u}_j^T \Sigma \mathbf{u}_j \\ \text{subject to:} \quad & \mathbf{u}_j^T \mathbf{u}_j = 1, \quad \text{for } j = 1, \dots, d,\end{aligned}\quad (9)$$

which is equivalent to

$$\max_{\mathbf{u}_1, \dots, \mathbf{u}_d, \lambda_1, \dots, \lambda_d} L = \sum_{j=1}^d \mathbf{u}_j^T \Sigma \mathbf{u}_j - \sum_{j=1}^d \lambda_j (\mathbf{u}_j^T \mathbf{u}_j - 1). \quad (10)$$

This quadratic optimization problem is solved by

$$\begin{aligned}\frac{\partial L}{\partial \mathbf{u}_j} &= \Sigma \mathbf{u}_j - \lambda_j \mathbf{u}_j = 0 \quad \text{for } j = 1, \dots, d. \\ \frac{\partial L}{\partial \lambda_j} &= 1 - \mathbf{u}_j^T \mathbf{u}_j = 0 \quad \text{for } j = 1, \dots, d.\end{aligned}\quad (11)$$

That is,  $\lambda_j$ 's are eigenvalues of  $\Sigma$  and  $\mathbf{u}_j$ 's are the corresponding eigenvectors.

Apparently, the objective function is maximized by selecting the  $d$  largest eigenvalues and retaining the corresponding eigenvectors.

## 3 An Example with Two Attributes

### 3.1 The Data

Let's consider a system with only two attributes (i.e.,  $N=2$ ). We have a sample data of size 1,000 (i.e.  $M=1,000$ ),  $\mathbf{X} \in \mathbb{R}^{1000 \times 2}$ .

The sample means of the two features are:

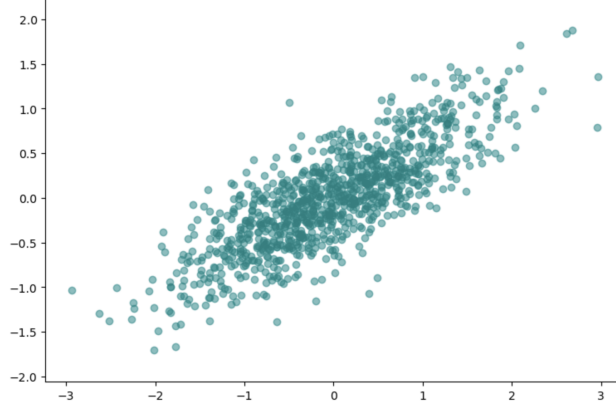
$$\bar{\mathbf{X}} = [\bar{x}_1, \bar{x}_2] = [-0.0366, -0.0047]. \quad (12)$$

The sample covariance matrix is

$$\Sigma = \frac{(\mathbf{X} - \bar{\mathbf{X}})^T (\mathbf{X} - \bar{\mathbf{X}})}{N - 1} = \begin{bmatrix} 0.7922 & 0.4049 \\ 0.4049 & 0.3207 \end{bmatrix}. \quad (13)$$

<sup>3</sup>Since  $\mathbf{P}$  is an orthogonal projection matrix,  $\mathbf{P} = \mathbf{P}^T$  and  $\mathbf{P}^2 = \mathbf{P}$ .  $\tilde{\mathbf{X}}\tilde{\mathbf{X}}^T$  is a constant.

<sup>4</sup> $\text{Tr}(\mathbf{A}\mathbf{B}) = \text{Tr}(\mathbf{B}\mathbf{A})$  if matrix multiplication  $\mathbf{A}\mathbf{B}$  and  $\mathbf{B}\mathbf{A}$  both exist.



**Figure 1:** A sample of two-feature data in size 1000

The correlation coefficient is

$$r = \frac{0.4049}{\sqrt{0.7922}\sqrt{0.3207}} = 0.803, \quad (14)$$

which means the two features are strongly correlated with each other.

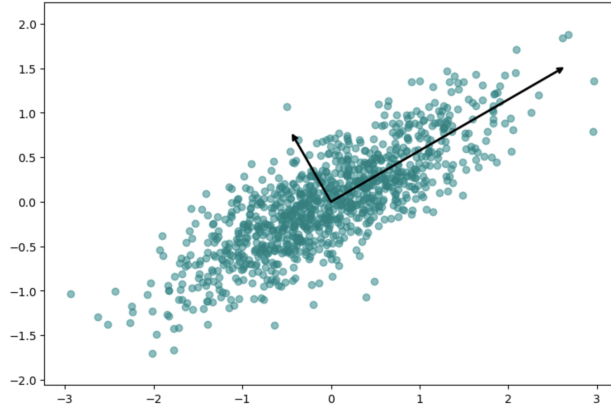
### 3.2 Eigen Decomposition of Covariance Matrix

We perform the eigen decomposition for  $\mathbf{X}$ 's covariance matrix. The eigenvectors are

$$\mathbf{V} = \begin{bmatrix} 0.8670 & -0.4984 \\ 0.4984 & 0.8670 \end{bmatrix} \approx \begin{bmatrix} \cos(30) & -\sin(30) \\ \sin(30) & \cos(30) \end{bmatrix}. \quad (15)$$

The corresponding eigenvalues in a diagonal matrix are

$$\mathbf{\Lambda} = \begin{bmatrix} 1.0249 & 0 \\ 0 & 0.0879 \end{bmatrix}. \quad (16)$$



**Figure 2:** A sample of two-feature data in size 1000

The eigenvectors define two perpendicular directions. As Figure 2 shows, the first direction,

$$\mathbf{v}_1 = \begin{bmatrix} 0.8670 \\ 0.4984 \end{bmatrix} \quad (17)$$

contains the most variation of the data. The corresponding eigenvalue  $\lambda_1 = 1.0249$  measures the amount of variance in this direction.

The second direction is defined by the second eigenvector,

$$\mathbf{v}_2 = \begin{bmatrix} -0.4984 \\ 0.8670 \end{bmatrix} \quad (18)$$

The variance in this direction is measured by the eigenvalue  $\lambda_2 = 0.0809$ .

### 3.3 PCA with Two Principal Components

Let's perform PCA on the zero-centered data  $\tilde{\mathbf{X}}$ . The two projection vectors, which are the eigenvectors of the covariance matrix of data  $\mathbf{X}$ , defines the projection matrix:

$$\mathbf{P} = [\mathbf{p}_1 \quad \mathbf{p}_2] = \begin{bmatrix} 0.8670 & -0.4984 \\ 0.4984 & 0.8670 \end{bmatrix}. \quad (19)$$

They respectively explain certain amount of variances, measured by their eigenvalues:

$$\mathbf{\Lambda} = \begin{bmatrix} \lambda_1 & 0 \\ 0 & \lambda_2 \end{bmatrix} = \begin{bmatrix} 1.0249 & 0 \\ 0 & 0.0879 \end{bmatrix}. \quad (20)$$

You must have found that

- the principal components are normalized linear combinations of  $\mathbf{X}$ 's features using the eigenvectors of  $\mathbf{X}$ 's sample covariance matrix
- the variances that principal components explain are the eigenvalues of  $\mathbf{X}$ 's sample covariance matrix.

The variance explained by the first principal component is:

$$\frac{1.0249}{1.0249 + 0.0879} = 92.1\%. \quad (21)$$

Therefore, using the first principal component as a feature seems also to be sufficient.

## 4 An Example of Feature Extraction using PCA

### 4.1 California Housing Dataset

Let's practice the feature extraction using PCA. We use the California Housing dataset provided by 'sklearn.dataset'.

```

1 .. _california_housing_dataset:
2
3 California Housing dataset
4 -----
5
6 **Data Set Characteristics:**
7
8 :Number of Instances: 20640
9
10 :Number of Attributes: 8 numeric, predictive attributes and the target
11
12 :Attribute Information:
13   - MedInc          median income in block group
14   - HouseAge        median house age in block group
15   - AveRooms         average number of rooms per household
16   - AveBedrms        average number of bedrooms per household
17   - Population       block group population
18   - AveOccup         average number of household members
19   - Latitude         block group latitude
20   - Longitude        block group longitude
21
22 :Missing Attribute Values: None
23
24 This dataset was obtained from the StatLib repository.
```

```

25 https://www.dcc.fc.up.pt/~ltorgo/Regression/cal\_housing.html
26
27 The target variable is the median house value for California districts, expressed in
    hundreds of thousands of dollars ($100,000).
28
29 This dataset was derived from the 1990 U.S. census, using one row per census
30 block group. A block group is the smallest geographical unit for which the U.S. Census
    Bureau publishes sample data (a block group typically has a population of 600 to 3,000
    people).
31
32 A household is a group of people residing within a home. Since the average number of rooms
    and bedrooms in this dataset are provided per household, these columns may take
    surprisingly large values for block groups with few households and many empty houses,
    such as vacation resorts.
33
34 It can be downloaded/loaded using the
35 :func:'sklearn.datasets.fetch_california_housing' function.
36
37 .. topic:: References
38
39 - Pace, R. Kelley and Ronald Barry, Sparse Spatial Autoregressions,
40   Statistics and Probability Letters, 33 (1997) 291-297

```

## 4.2 Data Normalization

Denote the data as  $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_N]$ . When we check the first few observations of the dataset, we found that those attributes have very different scales.

	MedInc	HouseAge	AveRooms	AveBedrms	Population	AveOccup	Latitude	Longitude
0	8.3252	41.0	6.984127	1.023810	322.0	2.555556	37.88	-122.23
1	8.3014	21.0	6.238137	0.971880	2401.0	2.109842	37.86	-122.22
2	7.2574	52.0	8.288136	1.073446	496.0	2.802260	37.85	-122.24
3	5.6431	52.0	5.817352	1.073059	558.0	2.547945	37.85	-122.25
4	3.8462	52.0	6.281853	1.081081	565.0	2.181467	37.85	-122.25

To eliminate the impact of varied scales of features, we normalize the data through centering (subtracting the sample mean) and scaling (dividing the sample standard deviation):

$$\widetilde{\mathbf{x}}_j = \frac{\mathbf{x}_j - \bar{x}_j}{s_j} \quad (22)$$

for  $j = 1, \dots, N$  before performing PCA. Here,

$$\bar{x}_j = \frac{\sum_{i=1}^M x_{i,j}}{M}. \quad (23)$$

is the sample mean of attribute  $j$ , and

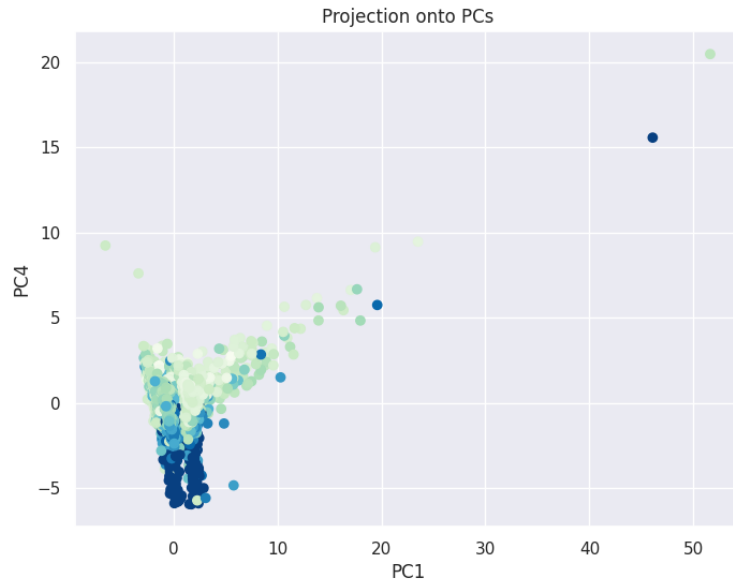
$$s_j = \sqrt{\frac{\sum_{i=1}^M (x_{i,j} - \bar{x}_j)^2}{M - 1}} \quad (24)$$

is the sample standard deviation.

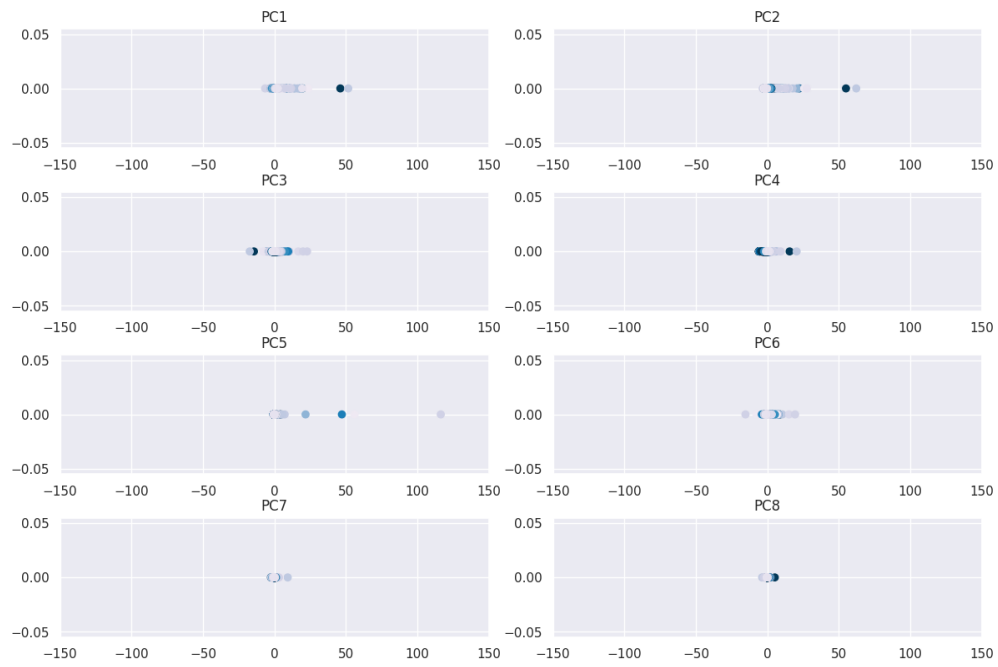
## 4.3 Data Visualization with Two PCs

We perform the PCA on the normalized data. Then, we visualize the data by showing the scatter plot of the first and fourth principal components, shown in Figure 3. The color of the data points is the median house value. The two PCs show the capability for interpreting the median house value.

We also visualize the principal components one by one below. The color of dots is corresponding to the house median price. It shows that the data have limited variation on the last two principal components.



**Figure 3:** Visualization of the data on the first two principal components



**Figure 4:** Visualization of the data on each principal component

#### 4.4 Interpret Principal Components

Let

$$[\mathbf{C}_1, \dots, \mathbf{C}_N] = \tilde{\mathbf{X}}\mathbf{P} = [\tilde{\mathbf{X}}\mathbf{p}_1, \dots, \tilde{\mathbf{X}}\mathbf{p}_N] \quad (25)$$

be the reconstruction of the normalized data. The  $j$ th principal component,  $\mathbf{C}_j$ , is the normalized linear combination of  $\tilde{\mathbf{X}}$ 's  $N$  attributes using the  $j$ th eigen vector of  $\mathbf{X}$ 's covariance matrix,  $\mathbf{p}_j$ :

$$\mathbf{C}_j = \tilde{\mathbf{X}}\mathbf{p}_j = [\tilde{\mathbf{x}}_1, \dots, \tilde{\mathbf{x}}_N]\mathbf{p}_j = [\tilde{\mathbf{x}}_1, \dots, \tilde{\mathbf{x}}_N] \begin{bmatrix} p_{1,j} \\ \vdots \\ p_{N,j} \end{bmatrix} = p_{1,j}\tilde{\mathbf{x}}_1 + \dots + p_{N,j}\tilde{\mathbf{x}}_N \quad (26)$$

for  $j = 1, \dots, N$ . Therefore, we can interpret  $\mathbf{p}_j$ 's as the loadings of  $\tilde{\mathbf{X}}$  for the linear transformation in eq. (25). Figure 5 visualizes  $\mathbf{p}_j$  for  $j = 1, \dots, 8$  as rows of the heat map. The larger the value of a cell, the higher the correlation between the corresponding principal component and feature. For example, the fifth principal component and the attribute AveOccup are strongly positively correlated. The fourth principal component and the feature MedInc are strongly negatively correlated.

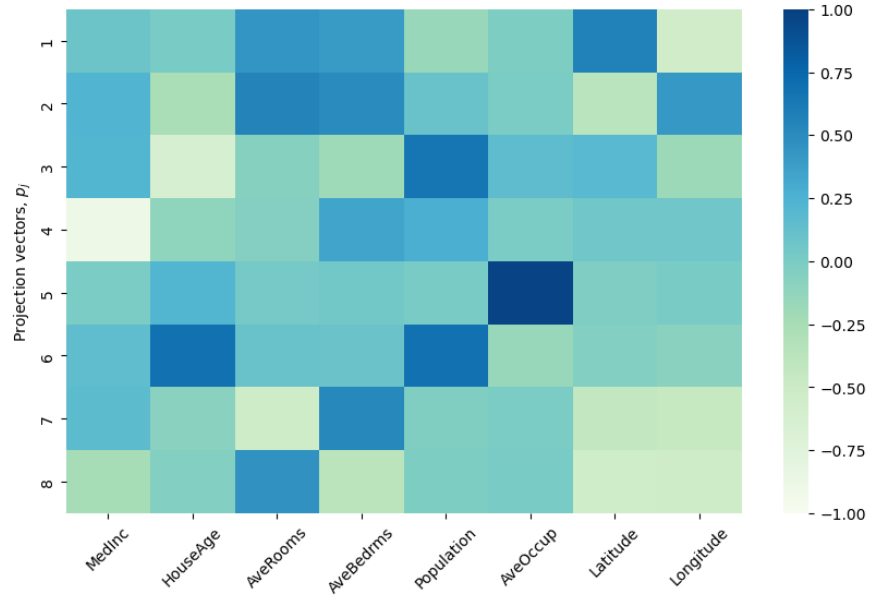
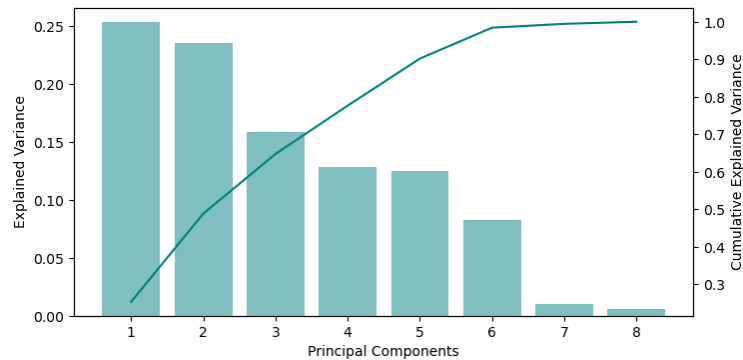


Figure 5: Interpret principal components



## 4.5 Explained Variance by Principal Components

In Figure 6, bars show the proportions of variance explained by individual principal components and the line shows the proportion of variances explained by the top  $d$  principal components. The figure shows the first principal component explains 25.34% of total variance. The top 5 principal components explain over 90% of total variance. Therefore, we might extract fewer principal components as new features for predicting the median value of houses. The principal components are independent because their correlation coefficients are zeros. If we build a multiple factor model using principal components for predicting the median house value, the collinearity issue is no longer an issue.

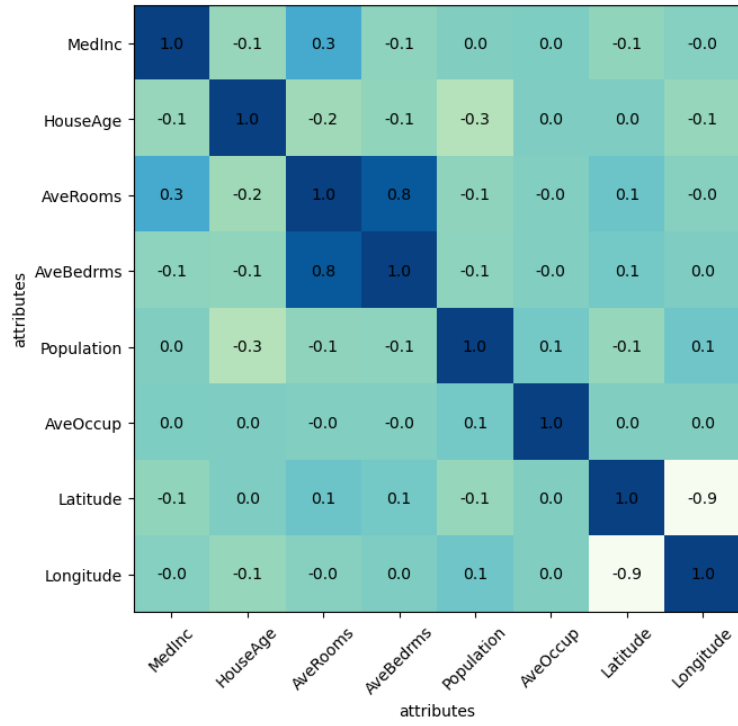


**Figure 6:** Explained variances by principal components

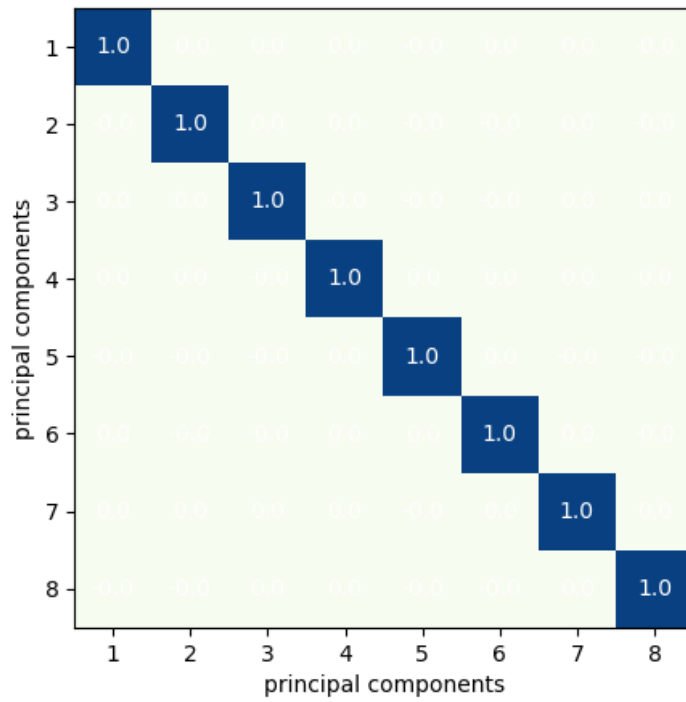
to be continued...

## References

- [1] Mehryar Mohri, Afshin Rostamizadeh, and Ameet Talwalkar. *Foundations of machine learning*. MIT press, 2018.



(a)



(b)

**Figure 7:** (a) correlation coefficients between attributes; (b) correlation coefficients between principal components