

### 3.7.4 PCA & regression

→ overview: Recall that a statistic is described as summarizing a random variable by mapping them to one value. Now consider summarizing a dataset w/ variables  $x_1, \dots, x_n$  into new w/ fewer variables. We can achieve this by performing principal component analysis (PCA). PCA is an unsupervised way to obtain new variables that summarize variability in a dataset. A more technical definition of PCA is: a procedure that finds a low-dimensional representation w/ a dataset w/o significant loss of information.

→ after explaining how PCA operates, we consider these new variables as predictors in a regression setting — that is called principal component regression (PCR).

→ for this lecture, we denote  $\theta_1, \dots, \theta_n$  as the latent variables in the original.

#### → PCA

→ just as a statistic summarizes RVs via a function, PCA summarizes variables via extraction as well. Specifically, a new variable  $z_i$  is created by taking linear combination of the original variables:

→ Denote such a new variable as  $z_{i,n}$ , which we call the  $n$ th principal component. Then  $z_{i,n}$  is defined as

$$z_{i,n} = \sum_{j=1}^n \theta_{i,j} x_{j,i}$$

→ The coefficients or this linear combination  $\theta_{i,1}, \dots, \theta_{i,n}$  are called the loadings of the  $n$ th principal component.

→ Then, for each observation  $i$  in the dataset, we can obtain an  $n$ th principal component score. To find the  $n$ th principal component score for the  $i$ th observation, evaluate the variables  $x_{j,i}$  at the corresponding values for the  $i$ th observation, i.e.

$$z_{i,n} = \sum_{j=1}^n \theta_{i,j} x_{j,i}$$

#### → First principal component

→ How do we determine the principal component loadings? The PCs, or principal components, are created in a sequence, so we start w/ the first principal component. The goal is for  $z_1$  to explain the largest portion of variability in a dataset, compared to subsequent ones. To achieve this, we maximize the (biased) sample variance of  $z_1$ .

→ In other words, the values of  $\theta_{1,1}, \dots, \theta_{1,n}$  are determined by maximizing

$$\frac{1}{n} \sum_{i=1}^n (z_{i,1} - \bar{z}_1)^2 \quad \text{c. } \bar{z}_1 = \text{all variables are centered}$$

$$= \frac{1}{n} \sum_{i=1}^n (\theta_{1,1} x_{1,i} + \dots + \theta_{1,n} x_{n,i})^2$$

→ However, we maximize while constraint by  $\sum_{i=1}^n \theta_{1,i}^2 = 1$  in order to obtain meaningful results.

→ To solve for these loadings, a technique known as eigen decomposition is often used. Learning this technique is not required for this space. However, note that the loadings are determined using the  $x_{1,1}, \dots, x_{1,n}$  variables only. In the absence of a response variable, PCA is an unsupervised learning (clustering) method.

#### → Second & subsequent principal components

→ Next consider the second principal component. The goal is for  $z_2$  to explain the next largest portion of remaining variability in a dataset that is not explained by  $z_1$ . As for the previous first principal component, this means:

→  $z_2$  is further constrained to be uncorrelated w/  $z_1$ . This is equivalent to saying that the vector of loadings for the first principal component is orthogonal or perpendicular to the vector of loadings for the second principal component, i.e.

$$\sum_{i=1}^n \theta_{1,i} \cdot \theta_{2,i} = 0 \Rightarrow z_1 \perp z_2$$

$$= \theta_1 \cdot \theta_2$$

$$\text{inner product}$$

→ The variability explained by the second principal component is less than the variability explained by the first principal component.

→ We can generalize this for all subsequent principal components:

→ All principal components are uncorrelated w/ one another. In other words,

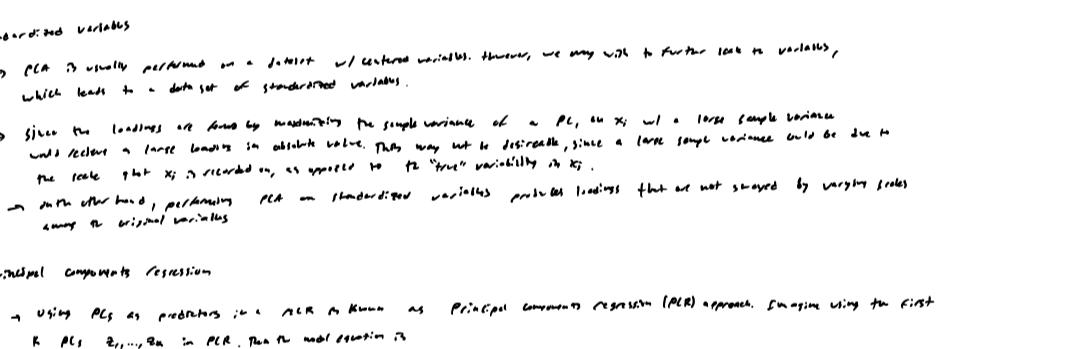
$$\sum_{i=1}^n \theta_{1,i} \cdot \theta_{k,i} = 0$$

for all pairs  $(k, 1, 2, \dots, n)$ .

→ The variability explained by each subsequent principal component is always less than the variability explained by its previous principal component.

→ The variability explained by each subsequent principal component — to some of original variables, all  $n$  principal components would collectively capture all of the variability from the dataset. The ideal situation is for the first few principal components to explain most of the dataset's variability. Thus, PCA is a dimension reduction technique as all  $n$  original variables are being summarized/compressed into fewer variables.

#### → Example



→ The first PC loadings are  $\theta_{1,1} = -0.4178$  &  $\theta_{1,2} = 0.8916$ . Therefore

$$z_1 = -0.4178 x_1 + 0.8916 x_2$$

→ The second PC loadings are  $\theta_{2,1} = -0.8916 x_1 + 0.4178 x_2$ . Hence,

$$z_2 = -0.8916 x_1 + 0.4178 x_2$$

→ In addition note that

→ The PCs represent the weight. That is a consequence of centering both variables before performing PCA.

→ The sum of squares of the (uncentered) first PC is 1 (because  $\theta_1 \cdot \theta_1 = 1$ ).

$$\theta_{1,1}^2 + \theta_{1,2}^2 = (-0.4178)^2 + (0.8916)^2 = 1$$

$$\theta_{1,1}^2 + \theta_{1,2}^2 = 1 \rightarrow 1 = 1$$

→ The PCs are uncorrelated. This means the dot product of the loading vectors is  $\theta_1 \cdot \theta_2 = 0$ :

$$\sum_{i=1}^n \theta_{1,i} \cdot \theta_{2,i} = \theta_{1,1} \cdot \theta_{2,1} + \dots + \theta_{1,n} \cdot \theta_{2,n}$$

↓ = 0 → 0

→ To better understand the PC scores, let's study one observation in the dataset: the 12th observation.

$$(x_{1,12}, x_{1,12}, x_{2,12}) = (-20.329, 18.164)$$

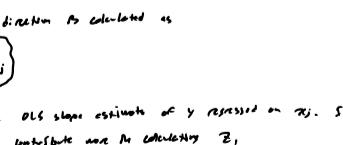
→ The observation has a first PC score of  $z_{1,12} = -0.4178(-20.329) + 0.8916(18.164)$

$$b = 9.19$$

→ A second PC score of  $z_{2,12} = -0.8916(-20.329) - 0.4178(18.164)$

$$b = 16.77$$

→ To grasp the meaning of these scores, we relate the coordinate system to first & second PCs. While we've learned to work in 2D space, remember that there are actually  $n$  dimensions. That results in a following coordinate where the second PC scores are plotted against the first PC:



→ We can see that a first PC score is the 2D-direction distance from an observation to the origin. On the other hand, a second PC score is the 2D-direction distance from an observation to the origin.

#### → Alternative Interpretations

→ We have viewed PCs as the directions where the data vary the most. However, as the observations live in  $n$ -dimensional space, we can also interpret PCs as the lower-dimensional surfaces that closest to the observations.

→ To demonstrate, imagine a dataset w/  $\geq 2$  of the centered variables. So, the obs live in 2D space. In this space, a  $n$ -dimensional surface can only be a line. That line that is closest to the obs is the one given by the first PC.

→ However, this surface can only be a line.

→ For any observation, the shortest distance to the line is the perpendicular distance to the line.

→ For any observation, the shortest distance to the line is the perpendicular distance to the line.

→ For any observation, the shortest distance to the line is the perpendicular distance to the line.

→ For any observation, the shortest distance to the line is the perpendicular distance to the line.

→ For any observation, the shortest distance to the line is the perpendicular distance to the line.

→ For any observation, the shortest distance to the line is the perpendicular distance to the line.

→ For any observation, the shortest distance to the line is the perpendicular distance to the line.

→ For any observation, the shortest distance to the line is the perpendicular distance to the line.

→ For any observation, the shortest distance to the line is the perpendicular distance to the line.

→ For any observation, the shortest distance to the line is the perpendicular distance to the line.

→ For any observation, the shortest distance to the line is the perpendicular distance to the line.

→ For any observation, the shortest distance to the line is the perpendicular distance to the line.

→ For any observation, the shortest distance to the line is the perpendicular distance to the line.

→ For any observation, the shortest distance to the line is the perpendicular distance to the line.

→ For any observation, the shortest distance to the line is the perpendicular distance to the line.

→ For any observation, the shortest distance to the line is the perpendicular distance to the line.

→ For any observation, the shortest distance to the line is the perpendicular distance to the line.

→ For any observation, the shortest distance to the line is the perpendicular distance to the line.

→ For any observation, the shortest distance to the line is the perpendicular distance to the line.

→ For any observation, the shortest distance to the line is the perpendicular distance to the line.

→ For any observation, the shortest distance to the line is the perpendicular distance to the line.

→ For any observation, the shortest distance to the line is the perpendicular distance to the line.

→ For any observation, the shortest distance to the line is the perpendicular distance to the line.

→ For any observation, the shortest distance to the line is the perpendicular distance to the line.

→ For any observation, the shortest distance to the line is the perpendicular distance to the line.

→ For any observation, the shortest distance to the line is the perpendicular distance to the line.

→ For any observation, the shortest distance to the line is the perpendicular distance to the line.

→ For any observation, the shortest distance to the line is the perpendicular distance to the line.

→ For any observation, the shortest distance to the line is the perpendicular distance to the line.

→ For any observation, the shortest distance to the line is the perpendicular distance to the line.

→ For any observation, the shortest distance to the line is the perpendicular distance to the line.

→ For any observation, the shortest distance to the line is the perpendicular distance to the line.

→ For any observation, the shortest distance to the line is the perpendicular distance to the line.

→ For any observation, the shortest distance to the line is the perpendicular distance to the line.

→ For any observation, the shortest distance to the line is the perpendicular distance to the line.

→ For any observation, the shortest distance to the line is the perpendicular distance to the line.

→ For any observation, the shortest distance to the line is the perpendicular distance to the line.

→ For any observation, the shortest distance to the line is the perpendicular distance to the line.

→ For any observation, the shortest distance to the line is the perpendicular distance to the line.

→ For any observation, the shortest distance to the line is the perpendicular distance to the line.

→ For any observation, the shortest distance to the line is the perpendicular distance to the line.

→ For any observation, the shortest distance to the line is the perpendicular distance to the line.

→ For any observation, the shortest distance to the line is the perpendicular distance to the line.

→ For any observation, the shortest distance to the line is the perpendicular distance to the line.

→ For any observation, the shortest distance to the line is the perpendicular distance to the line.

→ For any observation, the shortest distance to the line is the perpendicular distance to the line.

→ For any observation, the shortest distance to the line is the perpendicular distance to the line.

→ For any observation, the shortest distance to the line is the perpendicular distance to the line.

→ For any observation, the shortest distance to the line is the perpendicular distance to the line.

→ For any observation, the shortest distance to the line is the perpendicular distance to the line.

→ For any observation, the shortest distance to the line is the perpendicular distance to the line.

→ For any observation, the shortest distance to the line is the perpendicular distance to the line.

→ For any observation, the shortest distance to the line is the perpendicular distance to the line.

→ For any observation, the shortest distance to the line is the perpendicular distance to the line.

→ For any observation, the shortest distance to the line is the perpendicular distance to the line.

→ For any observation, the shortest distance to the line is the perpendicular distance to the line.

→ For any observation, the shortest distance to the line is the perpendicular distance to the line.

→ For any observation, the shortest distance to the line is the perpendicular distance to the line.

→ For any observation, the shortest distance to the line is the perpendicular distance to the line.

→ For any observation, the shortest distance to the line is the perpendicular distance to the line.

→ For any observation, the shortest distance to the line is the perpendicular distance to the line.

→ For any observation, the shortest distance to the line is the perpendicular distance to the line.

→ For any observation, the shortest distance to the line is the perpendicular distance to the line.

→ For any observation, the shortest distance to the line is the perpendicular distance to the line.

→ For any observation, the shortest distance to the line is the perpendicular distance to the line.

→ For any observation, the shortest distance to the line is the perpendicular distance to the line.