

PCA: Unsupervised Dimensionality Reduction

CS 530
Chapman
Spring 2021

Take  message for rest of course: PCA

- What are dimensionality reduction methods?
- What are these methods good for?
- How are they carried out?
- When to use and not to use them?

1

2

TABLE OF CONTENTS

- Review of Unsupervised Dimensionality Reduction Methods
- Principal Component Analysis (PCA)
 - Computing Principal Components
 - Scree Plot: How Many Components to Use?
 - Pros & Cons of PCA

PRINCIPAL COMPONENT ANALYSIS (PCA)

3

4

Dimensionality Reduction

If our data are composed of many features (dimensions), we often want to reduce the number of dimensions of our data. This often improves prediction accuracy, interpretability, or enables better visual inspection of the data.

Dimensionality reduction methods include:

- Feature Selection
- Principal Component Analysis (PCA)
- Linear Discriminant Analysis (LDA; supervised)
- Fischer's Linear Discriminant

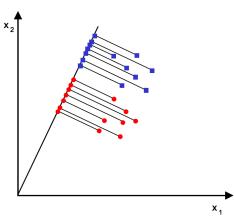


Image Source

Dimensionality Reduction

The term **feature selection**, as we saw for regression, refers to finding a subset of features— $x_{i_1}, x_{i_2}, \dots, x_{i_k}$ —that best explain the variability in the response variable, y . We could look at the p -values of the variables or run a subset selection method.

For supervised learning problems, we seek the best subset according to which the dataset can be learned.

Feature selection is different from dimensionality reduction.

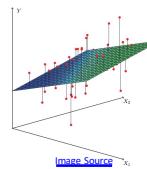


Image Source

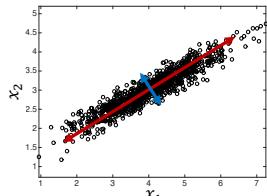
5

6

Dimensionality Reduction

Principal Component Analysis (PCA) finds *linear combinations of orthogonal variables* that explain the maximal amount of variance in the data.

In other words, PCA finds an alternative vector basis for the data, where the 1st dimension explains the maximal amount of variance in the data, the 2nd dimension explains the maximal amount of variance left to explain (orthogonal to the 1st dimension), and so on.

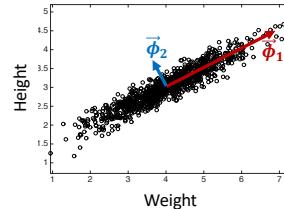


PCA: Explanatory Example

Say that you have a dataset with some population's height and weight.

The direction with most the variance in this data—the *1st principal component*—is designated in burgundy, $\vec{\phi}_1$.

The remaining variance—the *second principal component* in this 2D data—is then designated in blue, $\vec{\phi}_2$.



7

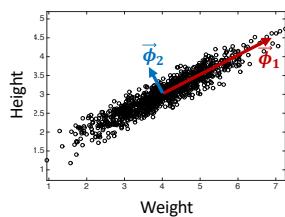
PCA: Explanatory Example

Here we have 1000 samples with 2 features (Height and Weight) each. So, we could describe these data as

$$\vec{x} = \begin{bmatrix} x_{1,1} & x_{1,2} \\ \vdots & \vdots \\ x_{1000,1} & x_{1000,2} \end{bmatrix}.$$

\vec{x} is thus *samples x features*, and generally $n \times p$:

$$\mathbb{R}^{n \times p} \ni \vec{x} = \begin{bmatrix} x_{1,1} & x_{1,2} & \cdots & x_{1,p} \\ x_{2,1} & x_{2,2} & \cdots & x_{2,p} \\ \vdots & \vdots & \ddots & \vdots \\ x_{n,1} & x_{n,2} & \cdots & x_{n,p} \end{bmatrix}$$



First Principal Component

For a given set of *features* (not samples), $\vec{x} = [\vec{x}_1, \vec{x}_2, \dots, \vec{x}_p]$ —so, each $\vec{x}_i = (x_{1,i}, x_{2,i}, \dots, x_{n,i})^T \in \mathbb{R}^n$ is a feature

$$\vec{x}_1, \vec{x}_2, \dots, \vec{x}_p \in \mathbb{R}^n$$

the **first principal component**, \vec{z}_1 , is the normalized linear combination of the features, $\vec{x}_1, \vec{x}_2, \dots, \vec{x}_p$, resulting in the maximal variance:

$$\mathbb{R}^n \ni \vec{z}_1 = \phi_{1,1}\vec{x}_1 + \phi_{2,1}\vec{x}_2 + \cdots + \phi_{p,1}\vec{x}_p = \vec{x} \cdot \vec{\phi}_1.$$

Here $\vec{\phi}_1 = (\phi_{1,1}, \phi_{2,1}, \dots, \phi_{p,1})^T \in \mathbb{R}^p$ are the *loadings*.

How do we choose \vec{z}_1 that (naively) maximizes variance?

9

First Principal Component

$\vec{z}_1 = \vec{x} \cdot \vec{\phi}_1$ is termed *normalized* linear combination because the coefficients, $\vec{\phi}_1$ —the *loadings* of the first principal component—describe a unit vector.

$$\sum_{j=1}^p \phi_{j,1}^2 = 1 \text{ or } \|\vec{\phi}_1\| = 1.$$

Remember, we want to maximize the variance of \vec{z}_1 .

Why do we need the constraint? In other words, what would happen if we do not put the constraint $\|\vec{\phi}_1\| = 1$?

10

First Principal Component

$\vec{z}_1 = \vec{x} \cdot \vec{\phi}_1$ is termed *normalized* linear combination because the coefficients, $\vec{\phi}_1$ —the *loadings* of the first principal component—describe a unit vector.

$$\sum_{j=1}^p \phi_{j,1}^2 = 1 \text{ or } \|\vec{\phi}_1\| = 1.$$

Without the constraint that $\|\vec{\phi}_1\| = 1$, setting $\vec{\phi}_1$ arbitrarily large would make $\vec{z}_1 = \vec{x} \cdot \vec{\phi}_1$ arbitrarily large too (for given data, \vec{x}) and thus \vec{z}_1 's variance would be arbitrarily large as well. So, without limiting $\|\vec{\phi}_1\|$, the solution to our optimization problem would be $\|\vec{\phi}_1\| \rightarrow \infty$, and so $\|\vec{z}_1\| \rightarrow \infty$. And that is not the solution we are seeking.

11

12

First Principal Component

Below we assume that \vec{x} was centered over each feature (column, dimension), and thus has a mean of 0.

If \vec{x} is centered over each feature, what can we do?

First Principal Component

Below we assume that \vec{x} was centered over each feature (column, dimension), and thus has a mean of 0. (if not, we subtract overall columns mean from all samples.)

We want to maximize the variance of the first principal component. We must solve the following constrained optimization problem:

$$\vec{\phi}_1 = \underset{\|\vec{\phi}_1\|=1}{\operatorname{argmax}} (\|\vec{x}^T \vec{\phi}_1\|^2) \quad (\text{i.e., subject to } \|\vec{\phi}_1\| = 1),$$

which is

$$\begin{aligned} \vec{\phi}_1 &= \underset{\|\vec{\phi}_1\|=1}{\operatorname{argmax}} \|\vec{x} \cdot \vec{\phi}_1\|^2 \\ &= \underset{\|\vec{\phi}_1\|=1}{\operatorname{argmax}} \left\{ \vec{\phi}_1^T \vec{x}^T \vec{x} \vec{\phi}_1 \right\} \underset{\|\vec{\phi}_1\|=1}{=} \underset{\|\vec{\phi}_1\|=1}{\operatorname{argmax}} \left\{ \frac{\vec{\phi}_1^T \vec{x}^T \vec{x} \vec{\phi}_1}{\vec{\phi}_1^T \vec{\phi}_1} \right\}. \end{aligned}$$

13

Reminder: Positive Semi-Definite Matrix

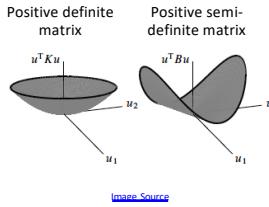
Definition:

A symmetric square matrix, $\vec{a} \in \mathbb{R}^{d \times d}$, is *positive definite* if

$$\vec{v}^T \vec{a} \vec{v} > 0 \text{ for all } \vec{v} \in \mathbb{R}^d \setminus \{0\}.$$

A symmetric square matrix, $\vec{a} \in \mathbb{R}^{d \times d}$, is *positive semi-definite* if

$$\vec{v}^T \vec{a} \vec{v} \geq 0 \text{ for all } \vec{v} \in \mathbb{R}^d.$$



15

Eigenvectors and Eigenvalues

Every square (non-pathological) matrix, \vec{A} , corresponds to a linear transformation. So, typically, it rotates and scales a vector, \vec{y} :

$$\begin{bmatrix} 2 & 1 \\ 1 & 2 \end{bmatrix} \begin{bmatrix} y_1 \\ y_2 \end{bmatrix} = \begin{bmatrix} 2y_1 + y_2 \\ y_1 + 2y_2 \end{bmatrix}.$$

But, for \vec{A} 's eigenvectors $(1,1)^T$ and $(1,-1)^T$:

$$\begin{bmatrix} 2 & 1 \\ 1 & 2 \end{bmatrix} \begin{bmatrix} 1 \\ 1 \end{bmatrix} = \begin{bmatrix} 3 \\ 3 \end{bmatrix} = 3 \cdot \begin{bmatrix} 1 \\ 1 \end{bmatrix}$$

and

$$\begin{bmatrix} 2 & 1 \\ 1 & 2 \end{bmatrix} \begin{bmatrix} 1 \\ -1 \end{bmatrix} = \begin{bmatrix} 3 \\ -3 \end{bmatrix} = 3 \cdot \begin{bmatrix} 1 \\ -1 \end{bmatrix}.$$

So, \vec{A} only scales each eigenvector \vec{v} by its corresponding eigenvalue λ :

$$\vec{A}\vec{v} = \lambda\vec{v}.$$

First Principal Component

$\vec{C} = \vec{x}^T \vec{x}$ is positive semi-definite, so $\frac{\vec{\phi}_1^T \vec{x}^T \vec{x} \vec{\phi}_1}{\vec{\phi}_1^T \vec{\phi}_1} = \frac{\vec{\phi}_1^T \vec{C} \vec{\phi}_1}{\vec{\phi}_1^T \vec{\phi}_1}$ is the Rayleigh quotient of \vec{C} , which reaches a maximum at \vec{C} 's maximal eigenvalue, λ_{max} , with \vec{v}_{max} being its corresponding eigenvector.

In other words, it can be shown that

$$\frac{\vec{\phi}_1^T \vec{C} \vec{\phi}_1}{\vec{\phi}_1^T \vec{\phi}_1} \leq \lambda_{max}.$$

And

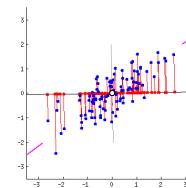
$$\frac{\vec{v}_{max}^T \vec{C} \vec{v}_{max}}{\vec{v}_{max}^T \vec{v}_{max}} = \lambda_{max}.$$

Where (again) λ_{max} and \vec{v}_{max} are the maximum eigenvalue of \vec{C} and its corresponding eigenvector

16

Visualizing the direction of 1st principal component

- What is the first principal component in the dataset below?



- The direction of maximal variance (magenta lines)

17

18

Further Principal Components

So $\vec{\phi}_1$ are the loadings of the 1st principal component. The m 'th principal component is calculated by finding the first ($m - 1$) principal component of matrix, and then computing

$$\vec{x}_{(m)} = \vec{x} - \sum_{l=1}^{m-1} \vec{x} \vec{\phi}_m \vec{\phi}_m^T,$$

in the same manner as we found the first principal component of \vec{x} :

$$\vec{\phi}_m = \underset{\|\vec{\phi}_m\|=1}{\text{argmax}} \left\| \vec{x}_{(m)} \cdot \vec{\phi}_m \right\|^2 = \underset{\|\vec{\phi}_m\|=1}{\text{argmax}} \left\{ \frac{\vec{\phi}_m^T \vec{x}_{(m)}^T \vec{x}_{(m)} \vec{\phi}_m}{\vec{\phi}_m^T \vec{\phi}_m} \right\}.$$

So, $\vec{\phi}_m$ are the loadings of the m 'th principal component.

It turns out that $\vec{\phi}_2, \dots, \vec{\phi}_p$ are simply the remaining eigenvectors of $\vec{x}^T \vec{x} = \vec{C}$.

All Principal Components

In matrix form, we write

$$\vec{z} = \vec{x} \vec{\phi},$$

where

$$\mathbb{R}^{p \times p} \ni \vec{\phi} = [\vec{\phi}_1, \vec{\phi}_2, \dots, \vec{\phi}_p],$$

and $\vec{\phi}_k$ is the k 'th eigenvector of $\vec{x}^T \vec{x} = \vec{C}$ (the covariance matrix of \vec{x}).

So, all together,

$$\vec{\phi}_1, \dots, \vec{\phi}_p$$

are the eigenvectors of $\vec{x}^T \vec{x}$ associated with eigenvalues

$$|\lambda_1| > |\lambda_2| > \dots > |\lambda_p|.$$

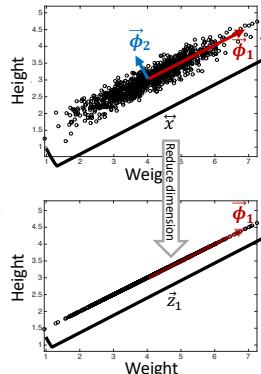
19

Principal Components

Geometrically, $\vec{\phi}_1$, which are the loadings of the first principal component, \vec{z}_1 , designate the direction along which the data, \vec{x} , vary the most.

\vec{z}_1 is then a reduced-dimension (1D) representation of the 2D \vec{x} data, along the most variable direction of \vec{x} .

For instance, $\vec{\phi}_1$ might be $(0.8 \text{ Weight}) + (0.6 \text{ Height})$. $0.8^2 + 0.6^2 = 1$, so the vector is properly normalized. If $\vec{\phi}_1 = (0.8, 0.6)$, $\vec{\phi}_2$ must be perpendicular to $\vec{\phi}_1$. So, $\vec{\phi}_2 = \pm(-0.6, 0.8)$.



20

Computing Principal Components

Solving

$$\vec{\phi}_1 = \underset{\|\vec{\phi}_1\|=1}{\text{argmax}} \|\vec{z}_1\|^2$$

can be done in two ways:

- Find the eigenvectors and eigenvalues of the covariance matrix, \vec{C} , of the data, \vec{x} .
- Use Singular Value Decomposition (SVD) on the data matrix, \vec{x} .

21

Computing Principal Components I

The covariance matrix, \vec{C} , is defined elementwise as

$$C_{ij} = \begin{cases} \text{variance of } x_i & i = j \\ \text{covariance of } x_i \text{ and } x_j & i \neq j \end{cases}$$

We can calculate \vec{C} using $\vec{C} = \frac{1}{n-1} \vec{x}^T \vec{x}$. The eigenvectors of \vec{C} are the directions of greatest variance in \vec{x} . Ordering these eigenvectors by the absolute values of their eigenvalues, we get the 1st, 2nd, 3rd, ... principal component directions. Diagonalizing \vec{C} (which is positive semidefinite):

$$\vec{C} = \vec{V} \vec{\Lambda} \vec{V}^T$$

$$\text{Then } \vec{V} = \begin{bmatrix} | & | & \cdots & | \\ v_1 & v_2 & \cdots & v_p \\ | & | & \cdots & | \end{bmatrix} \text{ and } \vec{\Lambda} = \begin{bmatrix} \lambda_1 & & 0 \\ & \lambda_2 & \\ 0 & & \ddots & \lambda_p \end{bmatrix},$$

with $|\lambda_1| > |\lambda_2| > \dots > |\lambda_p|$.

22

Computing Principal Components II

We can also compute the principal components directly using Singular Value Decomposition (SVD).

We can express any (sane...) matrix, including \vec{x} , as the product of 3 matrices:

$$\vec{x} = \vec{U} \cdot \vec{\Sigma} \cdot \vec{V}^T$$

Here \vec{U} and \vec{V} are unitary matrices—i.e., $\vec{U}^T = \vec{U}^{-1}$ and $\vec{V}^T = \vec{V}^{-1}$, so $\vec{V} \cdot \vec{V}^T = \vec{V} \cdot \vec{V}^{-1} = \vec{I}$ and the same for \vec{U} , hence their columns form orthonormal bases; and $\vec{\Sigma}$ is diagonal with the “singular values” of \vec{x} in descending order.

23

24

Computing Principal Components II

Plugging the SVD decomposition into the formula for \tilde{C} :

$$\begin{aligned}\tilde{C} &= \frac{1}{n-1} \tilde{x}^T \tilde{x} \\ &= \frac{1}{n-1} \tilde{V} \tilde{\Sigma} \tilde{U}^T \tilde{U} \tilde{\Sigma} \tilde{V}^T \\ &= \frac{1}{n-1} \tilde{V} \tilde{\Sigma}^2 \tilde{V}^T \\ &= \tilde{V} \left(\frac{\tilde{\Sigma}^2}{n-1} \right) \tilde{V}^T\end{aligned}$$

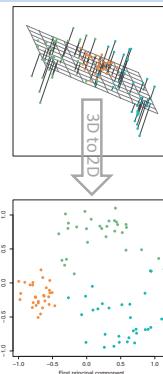
So, the columns of \tilde{V} are the principal directions, and

$$\frac{\tilde{\Sigma}_{i,i}^2}{n-1} = \lambda_i$$

ranks them in ascending order.

Geometrically, \vec{v}_1 corresponds to the direction in \tilde{x} with highest variance. \vec{v}_2 is the direction with highest variance remaining (perpendicular to \vec{v}_1); and so on.

Reducing Dimensions



So far, we saw how to represent the data, \tilde{x} , according to its “intrinsic directions”, as \tilde{z} .

If we want to reduce dimensionality from p to q , we take only the first (highest) q eigenvalues and their corresponding eigenvectors, creating truncated versions of \tilde{V} , $\tilde{\Sigma}$, and \tilde{U} : \tilde{V}^q , $\tilde{\Sigma}^q$, and \tilde{U}^q .

25

Reducing Dimensions

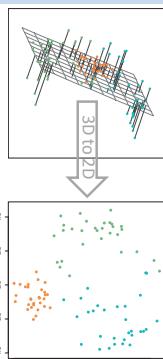
The reduced-dimension representation of \tilde{x} is \tilde{z}^q . \tilde{z}^q results in a linear combination of all the features for each ϕ_i (e.g., $\phi_1 = 0.8 \text{ Weight} + 0.6 \text{ Height}$) and not in feature selection.

Often, the first few dimensions represent much of the overall variance of the data (given overall p dimensions):

$$\frac{\sum_{i=1}^q \lambda_i}{\sum_{i=1}^p \lambda_i}.$$

In that case,

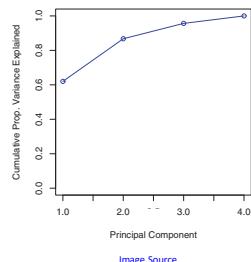
$$\tilde{x} \approx \tilde{z}^q \cdot \tilde{\Phi}^{q^{-1}}.$$



26

Scree Plot

But which approximation (i.e., which value of q) should we choose? Reducing to how many dimensions is good enough? We can decide how many principal components to use based on the cumulative proportion of variance explained, using a scree plot like the one on the right.



27

PCA Pros & Cons

Advantages of Principal Component Analysis:

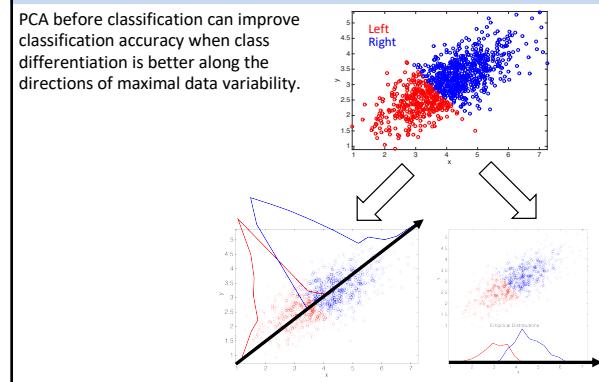
- Unsupervised preprocessing technique is labels agnostic, hence does not lead to double dipping (if used properly)
- Dimensionality reduction may lead to more interpretable data and data that is easier to visualize
- Running PCA does not rest on any special assumptions

Disadvantages of Principal Component Analysis:

- Principal components (linear combinations) can also sometimes be difficult to interpret directly
- Nonlinear structure can be hard to model with PCA (though nonlinear PCA methods with kernels and other extensions do exist)

PCA for Classification Preprocessing

PCA before classification can improve classification accuracy when class differentiation is better along the directions of maximal data variability.

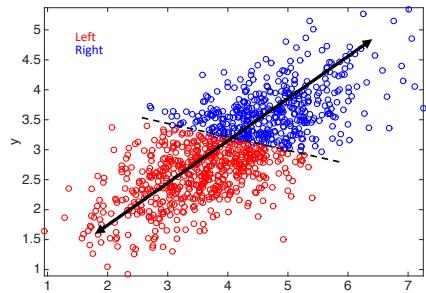


29

30

PCA for Classification Preprocessing

But PCA is not exposed to the labels. So, it will not do well when optimal class differentiation is not along the directions of maximal variability.



CONCLUSIONS

- Dimensionality reduction methods either remove some features of the data (feature selection) or create combinations of features and leave only ones optimal in some sense (e.g., PCA)
 - As a preprocessing step, this can potentially improve classification
 - It can also assist in visualizing the data
- Principal Component Analysis (PCA) is an unsupervised method (agnostic of labels) that finds “intrinsic” directions where the data’s variance is maximized along the 1st direction, the remaining variance is maximized along the 2nd direction, and so on

31

33

Take message for rest of course: PCA

- Dimensionality reduction methods attempt to simplify the data by reducing its number of features
- Preprocessing: they may reduce classification error
 - PCA: class differentiability along maximal variance direction
- PCA can improve data visualizability

34