## **One-dimensional PCA**

Archer Yang

McGill University

October 4, 2024

# Reading assignment: read chapter 10 of Mathematics for Machine Learning book.

#### **Dimension reduction**

- High-dimensional data is very common in real-life applications, e.g. in genomics, finance, e-commerce etc.
- Dimensions in high-dimensional data are often correlated it has an embedded lower-dimensional structure.
- Dimensionality reduction exploits such low dimensional structure
  - Compression
  - Visualization
  - Data generation
- We study principal component analysis (PCA) linear dimensionality reduction by Pearson (1901) and Hotelling (1933).
- One of the most commonly used techniques for compression and visualization.

#### **Dimension reduction**



(a) Dataset with  $x_1$  and  $x_2$  coordinates.

(b) Compressed dataset where only the  $x_1$  coordinate is relevant.

### **DERIVE PCA FROM FIRST PRINCIPLES**

#### **Problem Setting**

Consider an i.i.d. dataset

$$\mathcal{X} = \{\boldsymbol{x}_1, \cdots, \boldsymbol{x}_N\}$$

where 
$$\boldsymbol{x}_n = (x_{n1}, \cdots, x_{np})^\top \in \mathbb{R}^p$$
 for  $n = 1, \cdots, N$ .

For simplicity, we assume that the data is **centered** such that

$$\hat{\boldsymbol{\mu}} = rac{1}{N} \sum_{n=1}^{N} \boldsymbol{x}_n = \boldsymbol{0}$$

To achieve this, set

$$x_n \leftarrow x_n - \hat{\mu}$$

- Find projections  $\tilde{x}_n$  of data points  $x_n$ 
  - Similar to the original data points as possible
  - But have a significantly lower intrinsic dimensionality.

# PROJECTION PERSPECTIVE: ONE-DIMENSIONAL PROJECTION



#### Minimizing reconstruction error

- We'll start by looking for a one-dimensional projection.
- Project N points  $\boldsymbol{x}_n \in \mathbb{R}^p$ ,  $n = 1, \dots, N$  on to a line through the origin.
- We specify the line by a unit column vector

 $\boldsymbol{b}_1 = (b_{11}, \dots, b_{1p})^\top \in \mathbb{R}^p$ , with  $\|\boldsymbol{b}_1\| = 1$ , which forms a subspace U

$$U = \{z_1 \boldsymbol{b}_1 : \forall z_1 \in \mathbb{R}\}$$

The coordinate of the orthogonal projection of x<sub>n</sub> onto the line U (one-dimensional subspace spanned by b<sub>1</sub>) is

 $z_{1n} = \boldsymbol{b}_1^{\mathsf{T}} \boldsymbol{x}_n$ 

**The orthogonal projection**  $\tilde{x}_n$  is

$$\tilde{\boldsymbol{x}}_n \equiv \pi_U(\boldsymbol{x}_n) = \underbrace{\boldsymbol{b}}_1_{\text{direction of proj.}} \cdot \underbrace{\boldsymbol{b}}_1^\top \boldsymbol{x}_n_{z_{1n}:\text{coordinate of proj}}$$





#### Minimizing reconstruction error

If we approximate  $x_n$  using  $\tilde{x}_n$ , the squared approximation error is

$$\|\boldsymbol{x}_{n} - \tilde{\boldsymbol{x}}_{n}\|^{2} = \|\boldsymbol{x}_{n} - \pi_{U}(\boldsymbol{x}_{n})\|^{2} = \|\boldsymbol{x}_{n} - \boldsymbol{b}_{1}\boldsymbol{b}_{1}^{\mathsf{T}}\boldsymbol{x}_{n}\|^{2}$$

■ If average those  $||x_n - \tilde{x}_n||^2$  over all the points  $n = 1, \dots, N$ , we get the average squared approximation error, i.e. reconstruction error

$$J_1 \equiv \frac{1}{N} \sum_{n=1}^N \|\boldsymbol{x}_n - \tilde{\boldsymbol{x}}_n\|^2 = \frac{1}{N} \sum_{n=1}^N \|\boldsymbol{x}_n - \boldsymbol{b}_1 \boldsymbol{b}_1^\top \boldsymbol{x}_n\|^2 \quad \text{(reconstruction err.)}$$

How to choose a good subspace U in PCA? – find b<sub>1</sub> such that the reconstruction error is minimized

$$\min_{\boldsymbol{b}_1} \boldsymbol{J}_1 \Longleftrightarrow \min_{\boldsymbol{b}_1} \frac{1}{N} \sum_{n=1}^N \|\boldsymbol{x}_n - \boldsymbol{b}_1 \boldsymbol{b}_1^\top \boldsymbol{x}_n\|^2$$

#### Minimizing reconstruction error

 Here is an amazing result that bridges together two perspectives of PCA

$$\frac{1}{N}\sum_{n=1}^{N} \|\boldsymbol{x}_{n} - \boldsymbol{b}_{1}\boldsymbol{b}_{1}^{\mathsf{T}}\boldsymbol{x}_{n}\|^{2} = \frac{1}{N}\sum_{n=1}^{N} \|\boldsymbol{x}_{n}\|^{2} - \frac{1}{N}\sum_{n=1}^{N} \|\boldsymbol{b}_{1}^{\mathsf{T}}\boldsymbol{x}_{n}\|^{2}$$

which essentially claiming that

reconstruction err. const. projected variance  $\begin{bmatrix}
\frac{1}{N}\sum_{n=1}^{N} \|\boldsymbol{x}_{n} - \tilde{\boldsymbol{x}}_{n}\|^{2} = \begin{bmatrix}
\frac{1}{N}\sum_{n=1}^{N} \|\boldsymbol{x}_{n}\|^{2} - \begin{bmatrix}
\frac{1}{N}\sum_{n=1}^{N} z_{1n}^{2} \\
\frac{1}{N}\sum_{n=1}^{N} z_{1n}^{2}
\end{bmatrix}$ 

#### Proof

The squared approximation error for point  $x_n$  can be rewritten as

$$\begin{aligned} \|\boldsymbol{x}_{n} - \tilde{\boldsymbol{x}}_{n}\|^{2} &= \|\boldsymbol{x}_{n} - \boldsymbol{b}_{1}\boldsymbol{b}_{1}^{\mathsf{T}}\boldsymbol{x}_{n}\|^{2} \\ &= (\boldsymbol{x}_{n} - \boldsymbol{b}_{1}\boldsymbol{b}_{1}^{\mathsf{T}}\boldsymbol{x}_{n})^{\mathsf{T}}(\boldsymbol{x}_{n} - \boldsymbol{b}_{1}\boldsymbol{b}_{1}^{\mathsf{T}}\boldsymbol{x}_{n}) \\ &= \boldsymbol{x}_{n}^{\mathsf{T}}\boldsymbol{x}_{n} - \boldsymbol{x}_{n}^{\mathsf{T}}\boldsymbol{b}_{1}\boldsymbol{b}_{1}^{\mathsf{T}}\boldsymbol{x}_{n} - \boldsymbol{x}_{n}^{\mathsf{T}}\boldsymbol{b}_{1}\boldsymbol{b}_{1}^{\mathsf{T}}\boldsymbol{x}_{n} + \boldsymbol{x}_{n}^{\mathsf{T}}\boldsymbol{b}_{1}\boldsymbol{b}_{1}^{\mathsf{T}}\boldsymbol{b}_{1}\boldsymbol{b}_{1}^{\mathsf{T}}\boldsymbol{x}_{n} \\ &\stackrel{(i)}{=} \boldsymbol{x}_{n}^{\mathsf{T}}\boldsymbol{x}_{n} - \boldsymbol{x}_{n}^{\mathsf{T}}\boldsymbol{b}_{1}\boldsymbol{b}_{1}^{\mathsf{T}}\boldsymbol{x}_{n} \\ &= \|\boldsymbol{x}_{n}\|^{2} - \|\boldsymbol{b}_{1}^{\mathsf{T}}\boldsymbol{x}_{n}\|^{2} \\ &= \operatorname{const} - z_{1n}^{2} \end{aligned}$$

where (i) follows by  $\boldsymbol{x}_n^{\mathsf{T}} \boldsymbol{b}_1 \boldsymbol{b}_1^{\mathsf{T}} \boldsymbol{b}_1 \boldsymbol{b}_1^{\mathsf{T}} \boldsymbol{x}_n \stackrel{\boldsymbol{b}_1^{\mathsf{T}} \boldsymbol{b}_1 = 1}{=} \boldsymbol{x}_n^{\mathsf{T}} \boldsymbol{b}_1 \boldsymbol{b}_1^{\mathsf{T}} \boldsymbol{x}_n$  and (ii) follows by  $z_{1n} = \boldsymbol{b}_1^{\mathsf{T}} \boldsymbol{x}_n$ .

#### Sample variance of projected coordinates

• (Homework) Show that the term  $\frac{1}{N} \sum_{n=1}^{N} z_{1n}^2$  indeed is equal to the sample variance of the coordinates  $\{z_{11}, \dots, z_{1N}\}$  of N projections.

$$\frac{1}{N}\sum_{n=1}^{N}z_{1n}^{2} = \frac{1}{N}\sum_{n=1}^{N}\left(z_{1n} - \frac{1}{N}\sum_{n=1}^{N}z_{1n}\right)^{2} \equiv \widehat{\operatorname{Var}}(z_{11}, \dots, z_{1N})$$

**Hint:** recall that the data is centered; also note that we use factor  $\frac{1}{N}$  instead of  $\frac{1}{N-1}$  in the sample variance.

Note that Var stands for the sample variance, rather than the population variance Var.

(Homework) Show that we also have

$$\frac{1}{N}\sum_{n=1}^{N} \|\boldsymbol{x}_{n} - \tilde{\boldsymbol{x}}_{n}\|^{2} = \frac{1}{N}\sum_{n=1}^{N} \|\boldsymbol{x}_{n}\|^{2} - \frac{1}{N}\sum_{n=1}^{N} \|\tilde{\boldsymbol{x}}_{n}\|^{2}$$

i.e. norm of reconstruction is equal to norm of coordinate. (If you draw it, this is obvious).

#### Yin and Yang 阴阳 – data's duality

Minimizing the average squared reconstruction error is equivalent to maximizing the variance of the projection.

$$\min_{\boldsymbol{b}_1} J_1 \Longleftrightarrow \max_{\boldsymbol{b}_1} \widehat{\operatorname{Var}}(z_{11}, \cdots, z_{1N})$$

i.e.

$$\min_{\boldsymbol{b}_1} \frac{1}{N} \sum_{n=1}^N \|\boldsymbol{x}_n - \boldsymbol{b}_1 \boldsymbol{b}_1^\top \boldsymbol{x}_n\|^2 \Longleftrightarrow \max_{\boldsymbol{b}_1} \frac{1}{N} \sum_{n=1}^N \|\boldsymbol{b}_1^\top \boldsymbol{x}_n\|^2$$

Choosing a subspace to maximize the projected variance, or minimize the reconstruction error, is called principal component analysis (PCA).



**ChatGPT:** (experience PCA in a dual light) – one side radiates with the expansion of variance, while the other elegantly narrows down reconstruction errors. A visual symphony of data's duality.



**Archer:** (poor man's version?) – these figures show that maximizing the variance of the projection is equivalent to minimizing the average squared reconstruction error.

## **MAXIMUM VARIANCE PERSPECTIVE**

Accordingly, let's find the projection direction b<sub>1</sub> that maximizes the variance!

$$V_1 \equiv \widehat{\operatorname{Var}}(z_{11}, \dots, z_{1N}) = \frac{1}{N} \sum_{n=1}^N z_{1n}^2 = \frac{1}{N} \sum_{n=1}^N \| \boldsymbol{b}_1^{\mathsf{T}} \boldsymbol{x}_n \|^2$$

The corresponding optimization problem is

$$\max_{\boldsymbol{b}_1} \boldsymbol{V}_1 \equiv \frac{1}{N} \sum_{n=1}^N \|\boldsymbol{b}_1^{\mathsf{T}} \boldsymbol{x}_n\|^2 \tag{1}$$

For simplicity, rewrite the problem in a matrix format

$$V_{1} = \frac{1}{N} \sum_{n=1}^{N} \|\boldsymbol{b}_{1}^{\mathsf{T}}\boldsymbol{x}_{n}\|^{2}$$
  
$$= \frac{1}{N} \sum_{n=1}^{N} \boldsymbol{b}_{1}^{\mathsf{T}}\boldsymbol{x}_{n}\boldsymbol{x}_{n}^{\mathsf{T}}\boldsymbol{b}_{1} = \frac{1}{N} \sum_{n=1}^{N} \boldsymbol{b}_{1}^{\mathsf{T}}\boldsymbol{x}_{n}\boldsymbol{x}_{n}^{\mathsf{T}}\boldsymbol{b}_{1}$$
  
$$= \boldsymbol{b}_{1}^{\mathsf{T}} \left(\frac{1}{N} \sum_{n=1}^{N} \boldsymbol{x}_{n}\boldsymbol{x}_{n}^{\mathsf{T}}\right) \boldsymbol{b}_{1}$$
  
$$= \boldsymbol{b}_{1}^{\mathsf{T}} \left(\frac{1}{N} \boldsymbol{X}^{\mathsf{T}} \boldsymbol{X}\right) \boldsymbol{b}_{1}$$
  
$$= \boldsymbol{b}_{1}^{\mathsf{T}} S \boldsymbol{b}_{1}$$

#### Here S is the sample (data) covariance matrix

$$\boldsymbol{S} = rac{1}{N} \sum_{n=1}^{N} \boldsymbol{x}_n \boldsymbol{x}_n^{\mathsf{T}} = rac{1}{N} \boldsymbol{X}^{\mathsf{T}} \boldsymbol{X}$$

where

$$\boldsymbol{X} = \begin{bmatrix} \boldsymbol{x}_1^{\mathsf{T}} \\ \boldsymbol{x}_2^{\mathsf{T}} \\ \vdots \\ \boldsymbol{x}_N^{\mathsf{T}} \end{bmatrix} = \begin{bmatrix} x_{11} & x_{12} & \cdots & x_{1p} \\ x_{21} & x_{22} & \cdots & x_{2p} \\ \vdots & \vdots & \vdots \\ x_{N1} & x_{N2} & \cdots & x_{Np} \end{bmatrix}$$

Written in the matrix format, the optimization problem (1) becomes

$$\max_{\boldsymbol{b}_1} \boldsymbol{b}_1^\top \boldsymbol{S} \boldsymbol{b}_1$$

Note that arbitrarily increasing the magnitude of the vector b₁ increases b<sub>1</sub><sup>T</sup>Sb<sub>1</sub>. Thus we have to restrict b₁ to be a unit vector, i.e. ||b₁|| = 1

$$\max_{\boldsymbol{b}_1} \boldsymbol{b}_1^\top \boldsymbol{S} \boldsymbol{b}_1 \tag{2}$$

subject to  $\|\boldsymbol{b}_1\| = 1$ 

To solve problem (2), we introduce the Lagrangian

$$\mathcal{L}(\boldsymbol{b}_1, \lambda_1) = \boldsymbol{b}_1^{\mathsf{T}} \boldsymbol{S} \boldsymbol{b}_1 + \lambda_1 (1 - \boldsymbol{b}_1^{\mathsf{T}} \boldsymbol{b}_1)$$

The partial derivatives of L are

$$\frac{\partial \mathcal{L}}{\partial \boldsymbol{b}_1} = 2\boldsymbol{b}_1^{\mathsf{T}}\boldsymbol{S} - 2\lambda_1\boldsymbol{b}_1^{\mathsf{T}} \qquad \frac{\partial \mathcal{L}}{\partial \lambda_1} = 1 - \boldsymbol{b}_1^{\mathsf{T}}\boldsymbol{b}_1$$

Setting these partial derivatives to 0 gives

$$\boldsymbol{S}\boldsymbol{b}_1 = \lambda_1 \boldsymbol{b}_1 \qquad \boldsymbol{b}_1^{\mathsf{T}} \boldsymbol{b}_1 = 1$$

#### We see that

- **b**<sub>1</sub> is an **eigenvector** of the sample covariance matrix S.
- The Lagrange multiplier \(\lambda\_1\) is the corresponding eigenvalue, it is also the variance of the resulting projected coordinates

$$\boldsymbol{V}_1 = \boldsymbol{b}_1^{\top} \boldsymbol{S} \boldsymbol{b}_1 = \lambda_1 \boldsymbol{b}_1^{\top} \boldsymbol{b}_1 = \lambda_1$$

- Therefore, to maximize the variance V<sub>1</sub>, we choose the basis vector b<sub>1</sub> associated with the largest eigenvalue of the data covariance matrix S.
- The eigenvector **b**<sub>1</sub> is called the 1st principal component.

We can get the coordinate of the projection

$$z_{1n} = \boldsymbol{b}_1^\top \boldsymbol{x}_n \in \mathbb{R}$$

We can also get the approximation of x<sub>n</sub> by mapping the coordinate z<sub>1n</sub> back into data space

$$\tilde{\boldsymbol{x}}_n = \boldsymbol{b}_1 \boldsymbol{z}_{1n} = \boldsymbol{b}_1 \boldsymbol{b}_1^\top \boldsymbol{x}_n \in \mathbb{R}^p$$

which gives us the projected data point  $\tilde{x}_n$  in the original data space.

Remark: Although x
<sub>n</sub> is a p-dimensional vector, it only requires a single coordinate z<sub>1n</sub> to represent it with respect to the basis vector b<sub>1</sub> ∈ ℝ<sup>p</sup>.

# If we project the data to the first principal component b<sub>1</sub>, combining equation

$$J_1 \equiv \frac{1}{N} \sum_{n=1}^{N} \|\boldsymbol{x}_n - \tilde{\boldsymbol{x}}_n\|^2 = \frac{1}{N} \sum_{n=1}^{N} \|\boldsymbol{x}_n\|^2 - \frac{1}{N} \sum_{n=1}^{N} z_{1n}^2$$

and

$$\frac{1}{N}\sum_{n=1}^{N}z_{1n}^{2} = \frac{1}{N}\sum_{n=1}^{N}\|\boldsymbol{b}_{1}^{\top}\boldsymbol{x}_{n}\|^{2} = \boldsymbol{b}_{1}^{\top}\boldsymbol{S}\boldsymbol{b}_{1} = \lambda_{1}$$

we see that the corresponding reconstruction error is

$$J_{1} = \frac{1}{N} \sum_{n=1}^{N} ||\boldsymbol{x}_{n}||^{2} - \lambda_{1}$$

## **DATA CENTERING AND SCALING**

#### Data centering and scaling

- In PCA, the centering and scaling of data are important steps.
- But their necessity depends on the context and nature of your data.
- The original data matrix is

$$\boldsymbol{X} = \begin{bmatrix} x_{11} & x_{12} & \cdots & x_{1p} \\ x_{21} & x_{22} & \cdots & x_{2p} \\ \vdots & \vdots & & \vdots \\ x_{N1} & x_{N2} & \cdots & x_{Np} \end{bmatrix}$$

- Centering is almost always necessary in PCA.
- Involves substracting each entry by the corresponding column mean:

$$\begin{bmatrix} x_{11} - \bar{x}_1 & x_{12} - \bar{x}_2 & \cdots & x_{1p} - \bar{x}_p \\ x_{21} - \bar{x}_1 & x_{22} - \bar{x}_2 & \cdots & x_{2p} - \bar{x}_p \\ \vdots & \vdots & & \vdots \\ x_{N1} - \bar{x}_1 & x_{N2} - \bar{x}_2 & \cdots & x_{Np} - \bar{x}_p \end{bmatrix}$$

- The column mean of each column after centering becomes zero.
- Ensures that each variable contributes equally to the analysis so that PCA focuses on the variance of the data.
- Without, PCA might be influenced by variables that are on a larger scale.

#### Scaling - aka standardization

Involves dividing each entry by the corresponding column standard deviation:

$$\begin{bmatrix} \frac{x_{11} - \bar{x}_1}{\sqrt{\sqrt{\operatorname{Var}}(x_{11}, \dots, x_{N1})}} & \frac{x_{12} - \bar{x}_2}{\sqrt{\sqrt{\operatorname{Var}}(x_{12}, \dots, x_{N2})}} & \cdots & \frac{x_{1p} - \bar{x}_p}{\sqrt{\sqrt{\operatorname{Var}}(x_{1p}, \dots, x_{Np})} \\ \frac{x_{21} - \bar{x}_1}{\sqrt{\sqrt{\operatorname{Var}}(x_{11}, \dots, x_{N1})}} & \frac{x_{22} - \bar{x}_2}{\sqrt{\sqrt{\operatorname{Var}}(x_{12}, \dots, x_{N2})}} & \cdots & \frac{x_{2p} - \bar{x}_p}{\sqrt{\sqrt{\operatorname{Var}}(x_{1p}, \dots, x_{Np})} \\ \vdots & \vdots & & \vdots \\ \frac{x_{N1} - \bar{x}_1}{\sqrt{\sqrt{\operatorname{Var}}(x_{11}, \dots, x_{N1})}} & \frac{x_{N2} - \bar{x}_2}{\sqrt{\sqrt{\operatorname{Var}}(x_{12}, \dots, x_{N2})}} & \cdots & \frac{x_{Np} - \bar{x}_p}{\sqrt{\sqrt{\operatorname{Var}}(x_{1p}, \dots, x_{Np})} \end{bmatrix}$$

- The column variance of each column after centering becomes one.
- Crucial when the variables in data are on different scales (e.g., kg, km, F, C) or have different units of measurement.
- Scaling ensures that PCA gives equal weight to each variable, preventing variables with larger scales from dominating the PCA.

- Always center the data;
- Scale the data when variables are on different scales or units;
  - In imaging where all coordinates are in the same units, namely pixel intensities — there is no need to do such coordinate scaling.
- Be cautious with scaling for interpretability.

# WHY USE ORTHOGONAL PROJECTION IN RECONSTRUCTION?

#### Why orthogonal projection is optimal?

Recall that orthogonal projection is adopted to get reconstruction x<sub>n</sub>

$$J_{1} \equiv \frac{1}{N} \sum_{n=1}^{N} ||\boldsymbol{x}_{n} - \tilde{\boldsymbol{x}}_{n}||^{2} = \frac{1}{N} \sum_{n=1}^{N} ||\boldsymbol{x}_{n} - \boldsymbol{b}_{1} \boldsymbol{b}_{1}^{\mathsf{T}} \boldsymbol{x}_{n}||^{2}$$

But why use orthogonal projection?

in

- Consider an unknown linear projection method, the resulting coordinate of the projection for data point *x<sub>n</sub>* is *z<sub>1n</sub>*
- The corresponding reconstruction error is

$$J_1 = \frac{1}{N} \sum_{n=1}^{N} || \boldsymbol{x}_n - \boldsymbol{b}_1 z_{1n} ||^2$$

Assume  $b_1$  is given, find  $z_{1n}$  that minimize the reconstruction error

$$\min_{z_{11}, \dots, z_{1N}} \frac{1}{N} \sum_{n=1}^{N} \|\boldsymbol{x}_n - \boldsymbol{b}_1 z_{1n}\|^2$$

To solve this problem, compute the partial derivative and set it to zero

$$\frac{\partial J_1}{\partial z_{1n}} = -\frac{2}{N} (\boldsymbol{x}_n - \boldsymbol{b}_1 z_{1n})^{\mathsf{T}} \boldsymbol{b}_1 = -\frac{2}{N} (\boldsymbol{b}_1^{\mathsf{T}} \boldsymbol{x}_n - \boldsymbol{b}_1^{\mathsf{T}} \boldsymbol{b}_1 z_{1n}) = 0$$

Since  $\boldsymbol{b}_1^{\top} \boldsymbol{b}_1 = 1$ , the equation yields

$$z_{1n} = \boldsymbol{b}_1^\top \boldsymbol{x}_n$$

#### Consequently,

- The optimal linear projection  $\tilde{x}_n$  of  $x_n$  is an orthogonal projection.
- Optimal coordinates z<sub>1n</sub> of the projection x̃<sub>n</sub> are the coordinates of the orthogonal projection of the original data point x<sub>n</sub>.