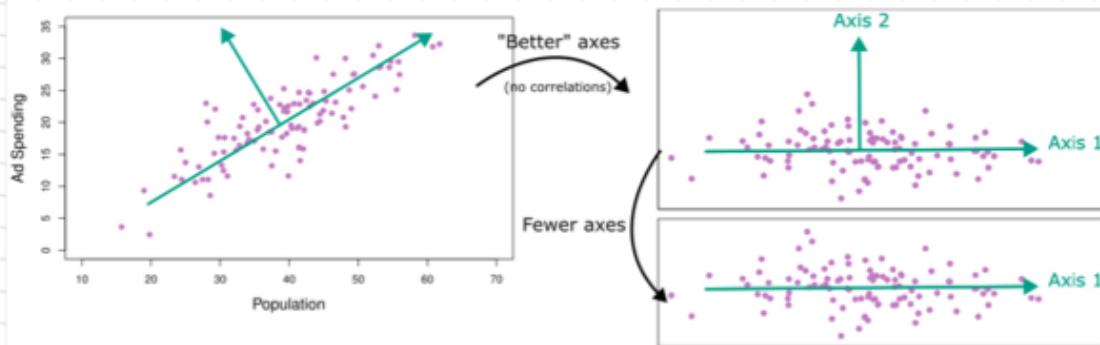


Lecture 10 - 01/10/24

Dimensionality Reduction & Visualization

- Normally our datasets are **HIGH DIMENSIONAL**



Principal Component Analysis

- Wish to represent data in a useful way for patterns
- Wanna see if data can be reduced to lower dimensions

PCA is an **UNSUPERVISED** method for such thing

- This means → **NO INFORMATION** about any class labels is used.
- PCA → finds a pattern
- PCA → looks for direction with **GREATEST VARIANCE**

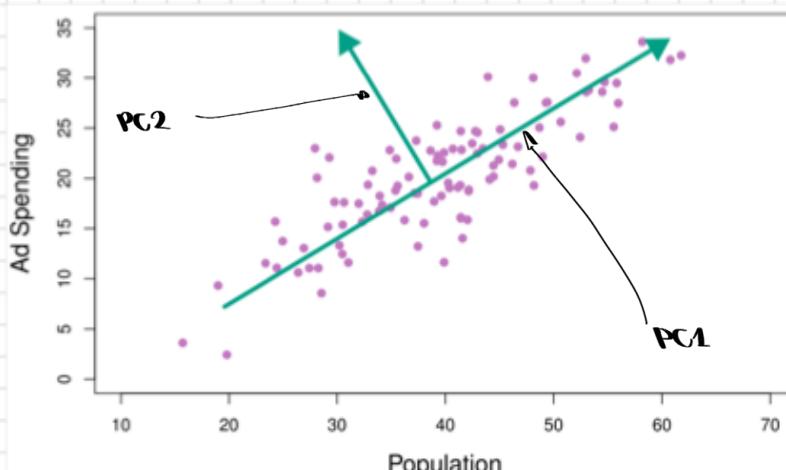
* FINDING NEW BASIS!

We want to build an **ORTHOGONAL BASIS** where new basis vectors are chosen to explain direction of the **GREATEST VARIANCE**.

* Project to Lower Dimension

First **K** principal components span a **K-dimensional subspace** that may be seen as the "best" **K-dimensional view** of data.

- PLA
- Consider a set of **p** features x_1, \dots, x_p each a real-valued random var.
 - Gives a **new set of p features**, principle components, each a **linear comb of original p**



Here the **PC1** score of data point i is:

$$z_{i1} = 0.839(\text{pop}_i - \bar{\text{pop}}) + 0.544(\text{ad}_i - \bar{\text{ad}}).$$

- It points in $(0.839, 0.544)$ in original
- New basis vector is **centered** at the column means for the data

The **PC2** is instead orthogonal:

$$z_{i2} = -0.544(\text{pop}_i - \bar{\text{pop}}) + 0.839(\text{ad}_i - \bar{\text{ad}})$$

First Principal Component

- Imagine that we have n data points, each with p features

- Data point i is:

$$\mathbf{x}_i = \begin{pmatrix} x_{i1} \\ \vdots \\ x_{ip} \end{pmatrix} \quad \xrightarrow{\text{SAMPLE MEAN}}$$

$$\bar{\mathbf{x}} = \sum_{i=1}^n \mathbf{x}_i / n$$

$$\mathbf{x}_i \cdot \mathbf{a}_1 = \|\mathbf{x}_i\| \|\mathbf{a}_1\| \cos \theta$$

- Wanna find direction of maximal variance

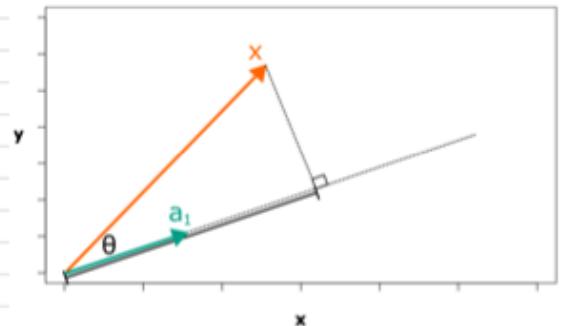
- We can use linear algebra to solve this

* Say that \mathbf{x}_i represents one data point

* Say that \mathbf{a}_1 is some other vector

* Do the dot Product

(to projects one vector onto another)

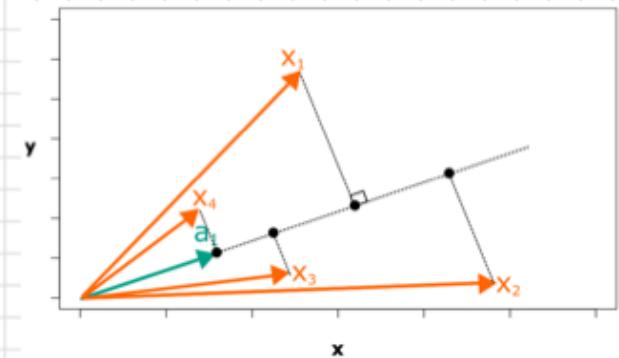


* We can now call \mathbf{a}_1 our first PC

$$\mathbf{a}_1 = \begin{pmatrix} a_{11} \\ \vdots \\ a_{1p} \end{pmatrix}$$

* If \mathbf{a}_1 points to max variance, $a_1 \cdot \mathbf{x}$ vary a lot

(to we wanna maximize $\text{Var}(a_1 \cdot \mathbf{x})$)



MAXIMIZATION

$$\text{Var}(\mathbf{a}_1^T \mathbf{X}) = \mathbf{a}_1^T S \mathbf{a}_1 \quad \text{WHERE} \quad S = \frac{1}{n} \sum_{i=1}^n (\mathbf{x}_i - \bar{\mathbf{x}})(\mathbf{x}_i - \bar{\mathbf{x}})^T = \text{Covariance matrix}$$

maximize this to constraint $\mathbf{a}_1^T \mathbf{a}_1 = 1$

we can use LAGRANGE MULTIPLIERS

$$F(\mathbf{a}_1, \lambda_1) = \mathbf{a}_1^T S \mathbf{a}_1 - \lambda_1 (\mathbf{a}_1^T \mathbf{a}_1 - 1)$$

function we want to maximize

take PARTIAL DERIVATIVE w.r.t \mathbf{a}_1 :

$$\frac{\partial F}{\partial \mathbf{a}_1} = 2S\mathbf{a}_1 - 2\lambda_1 \mathbf{a}_1 \xrightarrow{=0} S\mathbf{a}_1 - \lambda_1 \mathbf{a}_1 = 0$$

$$S\mathbf{a}_1 = \lambda_1 \mathbf{a}_1$$

EIGENVECTOR
CORRESPONDING
TO EIGENVALUE
 λ_1 IN DECOMP.

We also need to max F in λ_1 :

- Rewrite F using \mathbf{a}_1 satisfying $S\mathbf{a}_1 = \lambda_1 \mathbf{a}_1$

$$\begin{aligned} F(\mathbf{a}_1, \lambda_1) &= \mathbf{a}_1^T S \mathbf{a}_1 - \lambda_1 (\mathbf{a}_1^T \mathbf{a}_1 - 1) \\ &= \mathbf{a}_1^T \lambda_1 \mathbf{a}_1 - \lambda_1 \mathbf{a}_1^T \mathbf{a}_1 + \lambda_1 \\ &= \lambda_1. \end{aligned}$$

WE CHOOSE IT TO BE THE LARGEST EIGENVALUE

FINDING THE SECOND & SUBSEQUENT COMPONENT!

MAXIMIZE:

$$\text{Var}(\mathbf{a}_2^T \mathbf{x}) = \mathbf{a}_2^T \mathbf{S} \mathbf{a}_2 \quad \xrightarrow{\text{CONSTRAINTS}} \quad \mathbf{a}_2^T \mathbf{a}_2 = 1 \quad \& \quad \mathbf{a}_2^T \mathbf{a}_2 = 0$$

AS BEFORE USE LAGRANGE MULTIPLIERS:

$$F(\mathbf{a}_2, \lambda_2, \mu) = \mathbf{a}_2^T \mathbf{S} \mathbf{a}_2 - \lambda_2 (\mathbf{a}_2^T \mathbf{a}_2 - 1) - \mu (\mathbf{a}_2^T \mathbf{a}_2)$$

FROM HERE PARTIAL DERIVATIVES, MAX, ETC...

EIGENDECOMPOSITION

- With PCA we consider EIGENDECOMPOSITION OF \mathbf{S} :

$$\mathbf{S} = \mathbf{A} \Lambda \mathbf{A}^T$$

- We choose Λ to be the DIAGONAL MATRIX WITH EIGENVALUES $\lambda_1 \geq \dots \geq \lambda_p \geq 0$
 - The ORTHOGONAL $p \times p$ MATRIX $\mathbf{A} = [\mathbf{a}_1, \dots, \mathbf{a}_p]$ has EIGENVECTORS AS ITS COLUMNS
 - This MEANS THAT:
 - PCA CHOOSES EIGENDECOMPOSITION THAT ORDERS EIGENVECTORS ACCORDING TO DEC EIGENVAL
 - THE k^{th} PRINCIPAL COMPONENT IS:
- $$\mathbf{a}_k^T \mathbf{x} = a_{k1} x_1 + a_{k2} x_2 + \dots + a_{kp} x_p$$

- COEFFICIENTS OF PRINCIPAL COMPONENTS = LOADINGS.

- Vector of Loadings = EIGENVECTOR \mathbf{a}_k

- PRINCIPAL COMPONENT SCORES = OBS VAL OF $\mathbf{a}_k^T \mathbf{x}$

!!!!

- THE PCA ROTATES YOUR DATA AND PROJECTS IT TO $K \leq p$ DIMENSIONS
- PCA CREATES UNCORRELATED FEATURES THAT ARE LINEAR COMB OF EXISTING p .
- VARIANCE OF k^{th} PRINCIPAL COMPONENT IS THE k^{th} LARGEST EIGENVALUE $\rightarrow \text{Var}(\mathbf{a}_k^T \mathbf{x}) = \lambda_k$
- Each PC IS A VECTOR OF LEN p AS ORIGINAL DATA.
- IF USE ALL p PC'S CAN RECONSTRUCT DATA POINT \mathbf{x} :

$$\mathbf{x} = \sum_{i=1}^p (\mathbf{a}_i^T \mathbf{x}) \mathbf{a}_i \quad \xrightarrow{\text{i'm PC score for } \mathbf{x}, \text{ meaning coordinate in PC space}}$$

- DO THE SAME AS ABOVE SUST WITH FIRST m PC'S & WILL GET AN APPROXIMATION

Minimize approximation error (sum of squared error) by centering the data.



This corresponds to moving the coordinate system to the mean \bar{x} .



Mean and first four eigenvectors visualised:



Reconstructed images for increasing number of principal components



lets talk about Variance!

- total Sample Variance :

$$\text{Var}(x_1) + \dots + \text{Var}(x_p) = \text{Var}(z_1^T X) + \dots + \text{Var}(z_p^T X) = \lambda_1 + \dots + \lambda_p$$

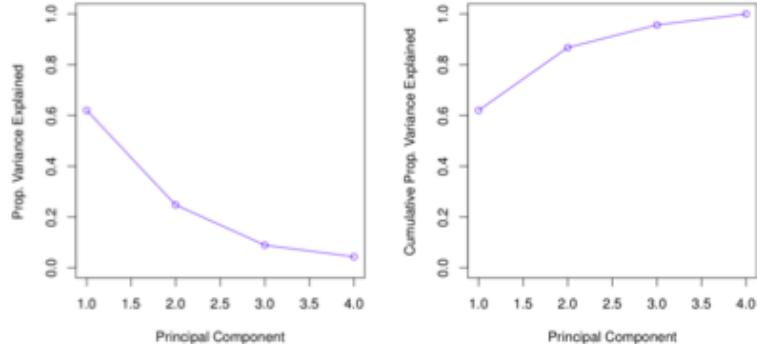
- Proportion of Variance explained by k^{th} PC :

$$\frac{\lambda_k}{\lambda_1 + \dots + \lambda_p}$$

NOTE

Total Var is the same for the original p's & p^{th} 's given that trace of $S \cdot P \Delta P^T$ is the same as trace of Δ

so How do we then decide # of components ??



Left: How much Var by each PC
Right: Cumulative of left-hand

Centering & Standardization

- A variable may be **centered** by :

$$\tilde{X} = X - \bar{X}$$

- Can then be scaled to unit variance :

$$\hat{X} = \frac{X}{\sqrt{\text{Var } X}} = (\text{Var } X)^{-\frac{1}{2}} X$$

- We can have the choice of **standardise features columnwise**:

- All features have unit var & correlation remains

- Another choice might be to **sphere the data**:

- All features have unit var & no correlation

