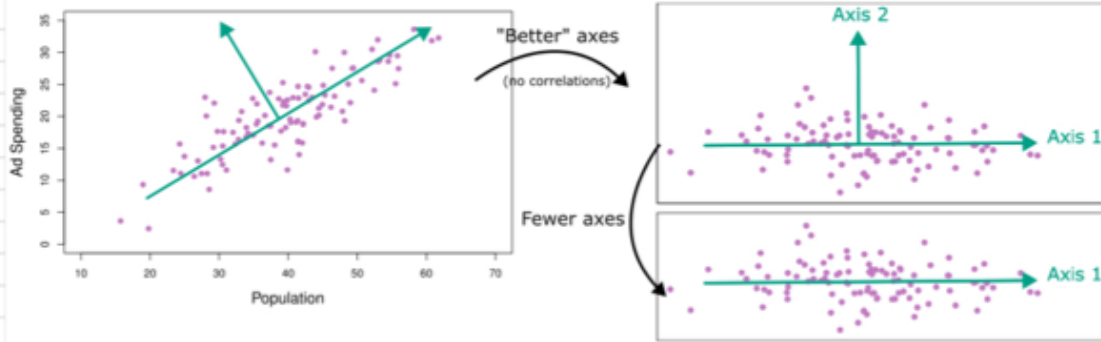


Lecture 10 - 01/10/24

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 \rightarrow **NO INFORMATION ABOUT ANY CLASS LABELS IS USED.**
- **PCA** \rightarrow **FINDS A PATTERN**
- **PCA** \rightarrow **LOOKS FOR DIRECTION WITH GREATEST VARIANCE**

* FINDING NEW BASIS!

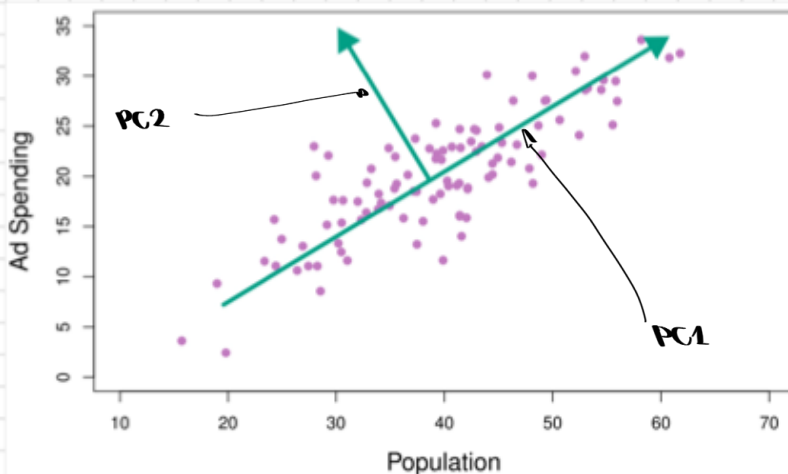
WE WANT TO **BUILD AN ORTHOGONAL BASIS** WHERE **NEW BASIS VECTORS** ARE CHOSEN TO EXPLAIN DIRECTION OF THE **GREATEST VARIANCE.**

* 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.**

- Consider a set of **p FEATURES** x_1, \dots, x_p EACH A **REAL-VALUED RANDOM VAR.**

PCA \rightarrow **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:

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

- 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:**

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

First Principal Component

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

- Data point i is:

$$x_i = \begin{pmatrix} x_{i,1} \\ \vdots \\ x_{i,p} \end{pmatrix} \xrightarrow{\text{SAMPLE MEAN}}$$

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

$$x_i \cdot a_1 = \|x_i\| \|a_1\| \cos \theta$$

- Wanna find direction of maximal variance

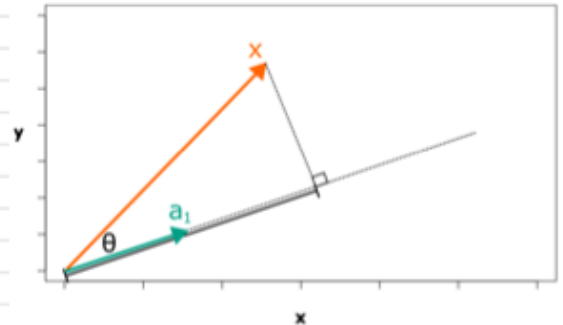
- We can use Linear Algebra to solve this

* Say that x_i represents one data point

* Say that a_1 is some other vector

* Do the dot product

↳ projects one vector onto another

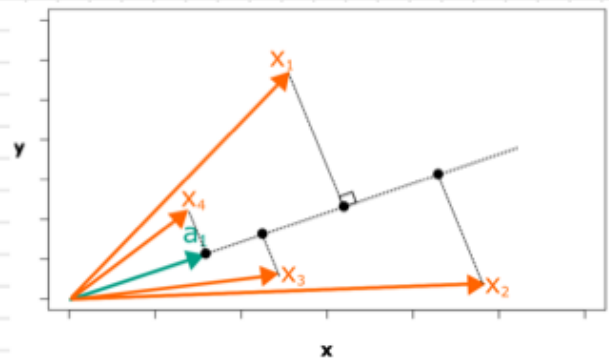


- We can now call a_1 our first PC

$$a_1 = \begin{pmatrix} a_{11} \\ \vdots \\ a_{1p} \end{pmatrix}$$

* If a_1 points to max variance, $a_1^T x$ vary a lot

↳ we wanna maximize $\text{Var}(a_1^T x)$



MAXIMIZATION

$$\text{Var}(a_1^T x) = a_1^T S a_1 \quad \text{where} \quad S = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})(x_i - \bar{x})^T = \text{Covariance matrix}$$

maximize this to constraint $a_1^T a_1 = 1$

we can use Lagrange multipliers

$$F(a_1, \lambda_1) = a_1^T S a_1 - \lambda_1 (a_1^T a_1 - 1)$$

↳ function we want to maximize

take partial derivative wrt a_1 :

$$\frac{\partial F}{\partial a_1} = 2S a_1 - 2\lambda_1 a_1 \xrightarrow{=0} S a_1 - \lambda_1 a_1 = 0$$

$$S a_1 = \lambda_1 a_1$$

EIGENVECTOR
CORRESPONDING
TO EIGENVALUE
 λ_1 in decomp.

We also need to max F in λ_1 :

- Rewrite F using a_1 satisfying $S a_1 = \lambda_1 a_1$

$$\begin{aligned} F(a_1, \lambda_1) &= a_1^T S a_1 - \lambda_1 (a_1^T a_1 - 1) \\ &= a_1^T \lambda_1 a_1 - \lambda_1 a_1^T a_1 + \lambda_1 \\ &= \lambda_1. \end{aligned}$$

↳ we choose it to be the largest eigenvalue

Finding the Second & subsequent component!

MAXIMIZE:

$$\text{Var}(\mathbf{a}_2^T \mathbf{X}) = \mathbf{a}_2^T \mathbf{S} \mathbf{a}_2 \xrightarrow{\text{constraints}} \mathbf{a}_2^T \mathbf{a}_2 = 1 \quad \& \quad \mathbf{a}_2^T \mathbf{a}_1 = 0$$

AS BEFORE USE LAGRANGE MULTIPLIERS:

$$F(\mathbf{a}_2, \lambda_2, \mu) = \mathbf{a}_2^T \mathbf{S} \mathbf{a}_2 - \lambda_2 (\mathbf{a}_2^T \mathbf{a}_2 - 1) - \mu (\mathbf{a}_2^T \mathbf{a}_1)$$

FROM HERE PARTIAL DERIVATIVES, MAX, ETC...

EIGENDECOMPOSITION

- WITH PCA WE CONSIDER EIGENDECOMPOSITION OF \mathbf{S} :

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

- WE CHOOSE \mathbf{A} 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 EIGVAL

- THE k^{th} PRINCIPAL COMPONENT IS:

$$\mathbf{a}_k^T \mathbf{X} = \mathbf{a}_{k1} X_1 + \mathbf{a}_{k2} X_2 + \dots + \mathbf{a}_{kp} X_p$$

- COEFFICIENTS OF PRINCIPLE 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} :

$$\mathbf{x} = \sum_{i=1}^p (\mathbf{a}_i^T \mathbf{x}) \mathbf{a}_i$$

i^{th} PC SCORE FOR \mathbf{x} , MEANING COORDINATE IN PC SPACE

- DO THE SAME AS ABOVE BUT 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 $\bar{\mathbf{x}}$.



Mean and first four eigenvectors visualised:



Reconstructed images for increasing number of principal components



lets talk about VARIANCE!

- TOTAL SAMPLE VARIANCE:

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

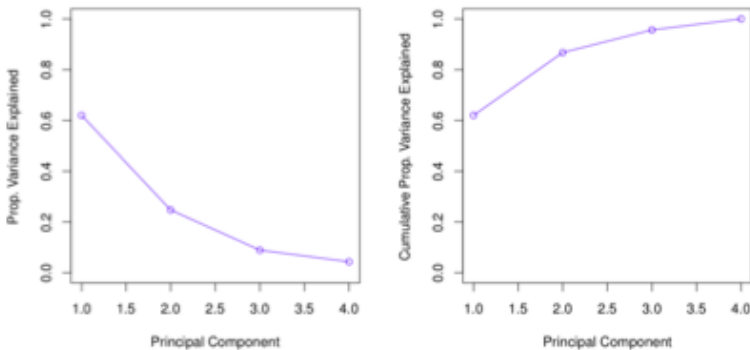
- PROPORTION OF VARIANCE EXPLAINED BY k^{th} PC:

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

NOTE

TOTAL VAR IS THE SAME FOR THE ORIGINAL p 's & p^{th} 's GIVEN THAT TRACE OF $S = P \Lambda P^T$ IS THE SAME AS TRACE OF Λ

60 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:

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

- CAN THEN BE SCALED TO UNIT VARIANCE:

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

- WE CAN HAVE THE CHOICE OF STANDARDISE FEATURES COMMONWISE:

- All features have UNIT VAR & CORRELATION REMAINS

- ANOTHER CHOICE MIGHT BE TO SPHERE THE DATA:

- All features have UNIT VAR & NO CORRELATION