

PRINCIPAL COMPONENT ANALYSIS

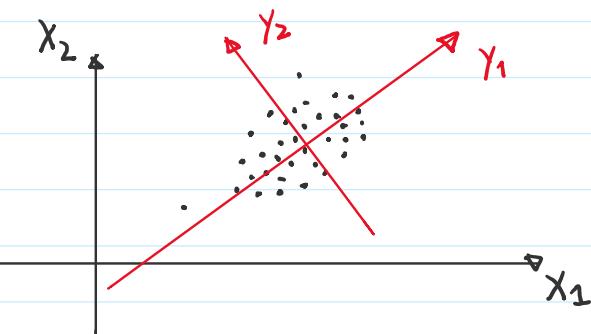
Sunday, April 7, 2019 11:34 AM

INTRODUCTION

THE PRINCIPAL COMPONENT ANALYSIS FOCUSES ON EXPLAINING THE COVARIANCE STRUCTURE OF A SET OF RANDOM VARIABLES (FROM WHICH OUR DATA ARE PRODUCED) THROUGH A FEW LINEAR COMBINATION OF THESE RANDOM VARIABLES.

ITS GENERAL OBJECTIVES ARE: 1. DATA DIMENSIONAL REDUCTION.
2. DATA INTERPRETATION. ANYWAY, THE CORE OF PCA IS: **FIND WHO CAUSES THE VARIANCE.**

SUPPOSE THAT WE HAVE A RANDOM VECTOR WITH 2 ELEMENTS $\underline{x}^t = [x_1, x_2]$, WHERE x_1 AND x_2 ARE RANDOM VARIABLES. FROM THIS RANDOM VECTOR WE CAN GENERATE A COLLECTION OF DATA, THAT IS, POINTS IN \mathbb{R}^2 . LOOK THE FIGURE BELOW, IN WHICH THE BLACK POINTS REPRESENT OUR DATA, IT SEEMS THAT THEY ARE FROM A MULTI-VARIATE GAUSSIAN DISTRIBUTION.



NOW, OUR GOAL IS: **FIND A DIRECTION (A VECTOR), ON WHICH THE VARIANCE OF PROJECTIONS OF DATA IS AS GREAT AS POSSIBLE; THEN REPEAT THE PROCEDURE ITERATIVELY, BUT THEY SHOULD BE MUTUALLY ORTHOGONAL.** IN OUR EXAMPLE, OBVIOUSLY y_1 AND y_2 , WHICH ARE DRAWN AS RED ARROWS IN THE FIGURE, ARE THE RIGHT ANSWER.

MATHEMATICAL FORMULATION

MATHEMATICAL FORMULATION

SUPPOSE THAT WE HAVE A RANDOM VECTOR

$$\underline{X}' = [X_1, X_2, \dots, X_p]$$

IN WHICH X_1, \dots, X_p ARE REAL-VALUED RANDOM VARIABLES. MOREOVER, \underline{X}' IS WITH COVARIANCE MATRIX Σ , THAT IS

$$\Sigma := \begin{bmatrix} \text{cov}(X_1, X_1) & \dots & \text{cov}(X_1, X_p) \\ \vdots & \ddots & \vdots \\ \text{cov}(X_p, X_1) & \dots & \text{cov}(X_p, X_p) \end{bmatrix}$$

THEN OUR PROBLEM CAN BE REWRITTEN AS AN OPTIMIZATION PROBLEM AS THE FOLLOWING:

P.C.A. PROBLEM

FIND $\underline{a}_1, \dots, \underline{a}_p \in \mathbb{R}^p$ WITH $\|\underline{a}_i\| = 1$, $i \in \{1, \dots, p\}$, SUCH THAT:

$$1) \underline{a}_1 = \underset{\underline{a} \in \mathbb{R}^p, \|\underline{a}\|=1}{\operatorname{argmax}} \text{Var}(\underline{a}' \underline{X}) \quad \text{EXPLICITY,}$$

$\underline{a}' \underline{X} = a_1 X_1 + a_2 X_2 + \dots + a_p X_p$
WHICH IS EXACTLY A LINEAR COMBINATION OF FOWER R.V.S

$$2) \underline{a}_2 = \underset{\underline{a} \in \mathbb{R}^p, \|\underline{a}\|=1}{\operatorname{argmax}} \text{Var}(\underline{a}' \underline{X})$$

$$\text{cov}(\underline{a}' \underline{X}, \underline{a}' \underline{X}) = 0$$

:

$$3) \underline{a}_p = \underset{\underline{a} \in \mathbb{R}^p, \|\underline{a}\|=1}{\operatorname{argmax}} \text{Var}(\underline{a}' \underline{X})$$

$$\text{cov}(\underline{a}' \underline{X}, \underline{a}' \underline{X}) = 0, \forall k = 1, \dots, p-1$$

IN FACT THE LAST PROBLEM IS NOT NECESSARY, BECAUSE THE SPACE IS WITH DIMENSION p . WHEN YOU HAVE CALCULATED THE FIRST $p-1$ DIRECTIONS, THE LAST ONE NATURALLY APPEARS.

SOLUTION OF PROBLEM

FOR SOLVING THIS OPTIMIZATION PROBLEM, WE NEED SOME LEMMAS:

LEMMA 1 THE COVARIANCE MATRIX Σ OF A RANDOM VECTOR $\underline{X}' = [X_1, X_2, \dots, X_p]$ IS POSITIVE SEMI-DEFINITE.

PROOF LET US USE μ TO DENOTE THE EXPECTED VALUES OF THE RANDOM VECTOR, THAT IS

$$\underline{\mu}' = (\mu_1, \mu_2, \dots, \mu_p)$$

$$\begin{aligned}\mu' &= (\mu_1, \mu_2, \dots, \mu_p) \\ &:= (\mathbb{E}[X_1], \mathbb{E}[X_2], \dots, \mathbb{E}[X_p])\end{aligned}$$

FOR $y \in \mathbb{R}^p$, WE HAVE

$$\begin{aligned}y' \Sigma y &= y' \mathbb{E}[(x - \mu)(x - \mu)'] y \\ &= \mathbb{E}[y'(x - \mu)(x - \mu)' y] \\ &= \mathbb{E}[(x - \mu)' y] \cdot (x - \mu)' y \quad \text{NOTICE } (x - \mu)' y \in \mathbb{R} \\ &= \mathbb{E}[(x - \mu)' y]^2 \geq 0\end{aligned}$$

BECAUSE y IS ARBITRARY, AND Σ IS SYMMETRIC, WE FINISH THE PROOF. ■

LEMMA 2 LET B BE A $p \times p$ POSITIVE SEMI-DEFINITIVE MATRIX WITH EIGENVALUES $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_p \geq 0$ AND ASSOCIATED NORMALIZED EIGENVECTORS e_1, e_2, \dots, e_p . THEN

$$(1) \quad \max_{\substack{x \neq 0 \\ x \perp x}} \frac{x'Bx}{x'x} = \lambda_1 \quad (\text{OBTAINED WHEN } x = e_1)$$

$$(2) \quad \min_{\substack{x \neq 0 \\ x \perp x}} \frac{x'Bx}{x'x} = \lambda_p \quad (\text{OBTAINED WHEN } x = e_p)$$

MOREOVER,

$$(3) \quad \max_{\substack{x \perp e_1, \dots, e_k \\ x \perp x}} \frac{x'Bx}{x'x} = \lambda_{k+1} \quad (\text{OBTAINED WHEN } x = e_{k+1}, \text{ WHERE } k = 1, 2, \dots, p-1)$$

PROOF THANKS TO **REAL SPECTRAL THEOREM**, WE CAN WRITE Σ AS:

$$B = P \Lambda P'$$

$$\text{WHERE } \Lambda = \begin{bmatrix} \lambda_1 & & \\ & \lambda_2 & \\ & & \ddots \\ & & & \lambda_p \end{bmatrix} \text{ AND } P = [e_1, e_2, \dots, e_p]$$

$$\text{IN PARTICULAR } P P' = P' P = I.$$

$$\text{DENOTE } B^{1/2} := P \Lambda^{1/2} P', \text{ WHERE } \Lambda^{1/2} := \begin{bmatrix} \sqrt{\lambda_1} & & \\ & \sqrt{\lambda_2} & \\ & & \ddots \\ & & & \sqrt{\lambda_p} \end{bmatrix},$$

ALSO $y := P' x$; IF $x \neq 0$, OBVIOUSLY $y \neq 0$. Thus

$$\frac{x'Bx}{x'x} = \frac{x' B^{1/2} B^{1/2} x}{x' x} = \frac{x' P \Lambda^{1/2} P' P \Lambda^{1/2} P x}{x' x} = \frac{y' \Lambda y}{y' y} = \dots$$

$$\frac{\underline{x}' \underline{x}}{\underline{x}' \underline{x}} = \frac{\underline{x}' \underline{x} \cdot \underline{x}}{\underline{x}' \underline{x}} = \frac{\underline{x}' \underline{x} \cdot \underline{x}}{\underline{x}' \underline{x}} = \frac{\underline{x}' \underline{x} \cdot \underline{x}}{\underline{x}' \underline{x}} = \frac{\underline{x}' \underline{x}}{\underline{x}' \underline{x}} = \dots$$

$$= \frac{\sum_{i=1}^p \lambda_i y_i^2}{\sum_{i=1}^p y_i^2} \leq \lambda_i \frac{\sum_{i=1}^p y_i^2}{y_i^2} = \lambda_i$$

SETTING $\underline{x} = \underline{e}_1$ GIVES

$$\underline{y} = P' \underline{e}_1 = \begin{bmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix}$$

THANKS A BUNCH
TO REAL SPECTRAL
THEOREM

UNDER THIS CONDITION, WE GET $\frac{\underline{y}' \underline{A} \underline{y}}{\underline{y}' \underline{y}} = \lambda_1$, WHICH
COMPLETES THE PROOF OF ①.

BY THE SIMILAR METHOD, WE CAN PROVE ②.

FOR ③, OBSERVE THAT $\underline{y} = P' \underline{x}$, THEN

$$\underline{x} = (P P') \underline{x} = P \underline{y} = y_1 \underline{e}_1 + y_2 \underline{e}_2 + \dots + y_p \underline{e}_p$$

SO $\underline{x} \perp \underline{e}_2, \dots, \underline{e}_p$ IMPLIES

$$0 = \underline{e}_i' \underline{x} = y_1 \underline{e}_i' \underline{e}_1 + y_2 \underline{e}_i' \underline{e}_2 + \dots + y_p \underline{e}_i' \underline{e}_p = y_i$$

THEN WE KNOW $y_i = 0, \forall i \leq p$. THEREFORE WE HAVE

$$\frac{\underline{x}' B \underline{x}}{\underline{x}' \underline{x}} = \frac{\sum_{i=1}^p \lambda_i y_i^2}{\sum_{i=1}^p y_i^2} = \frac{\sum_{i=p+1}^p \lambda_i y_i^2}{\sum_{i=p+1}^p y_i^2} \leq \lambda_{p+1}$$

TAKING $y_{p+1} = 1, y_{p+2} = \dots = y_p = 0$ GIVES THE MAXIMUM. ■

NOW, WE CAN FORMALLY STATE OUR FINAL RESULT AS

RESULT 3 THE SOLUTION OF P.C.A. PROBLEM IS

$$(\underline{a}_1, \underline{a}_2, \dots, \underline{a}_p) = (\underline{e}_1, \underline{e}_2, \dots, \underline{e}_p)$$

WHICH IS EQUIVALENT TO SAY:

THE i -TH PRINCIPAL COMPONENT OF \underline{x} IS $y_i := \underline{e}_i' \underline{x}$.

PROOF FROM LEMMA 1 WE KNOW $\underline{\Sigma}$ IS POSITIVE SEMI-DEFINITE;
AND $\max_{\underline{a} \in \mathbb{R}^p} \text{Var}(\underline{a}' \underline{x}) = \max_{\underline{a} \in \mathbb{R}^p} \underline{a}' \underline{\Sigma} \underline{a} = \max_{\underline{a} \in \mathbb{R}^p} \frac{\underline{a}' \underline{\Sigma} \underline{a}}{\|\underline{a}\|^2}$

$$\text{AND } \max_{\substack{\underline{a} \in \mathbb{R}^P \\ \|\underline{a}\|=1}} \text{Var}(\underline{a}' \underline{X}) = \max_{\substack{\underline{a} \in \mathbb{R}^P \\ \|\underline{a}\|=1}} \underline{a}' \underline{\Sigma} \underline{a} = \max_{\substack{\underline{a} \in \mathbb{R}^P \\ \|\underline{a}\|=1}} \frac{\underline{a}' \underline{\Sigma} \underline{a}}{\underline{a}' \underline{a}} \xrightarrow{\substack{\text{COVARIANCE} \\ \text{UNDER LINEAR} \\ \text{TRANSFORM}}} \text{NORMALIZATION}$$

THEN, IF WE CAN SHOW THE EQUIVALENT RELATIONSHIP:

$$\text{Cov}(Y_i, Y_k) = 0 \Leftrightarrow \underline{e}_i \perp \underline{e}_k$$

THEN WE CAN DIRECTLY APPLY LEMMA 2 TO GET THE RESULT; LUCKILY,

$$\text{Cov}(Y_i, Y_k) = \underline{e}_i' \underline{\Sigma} \underline{e}_k = \underline{e}_i' \lambda_k \underline{e}_k = \lambda_k \underline{e}_i' \underline{e}_k = 0.$$

WHICH COMPLETES THE PROOF. ■

CHARACTERIZATION OF TOTAL VARIANCE

HERE WE GIVE ALSO AN EASY BUT IMPORTANT RESULT ABOUT THE CHARACTERIZATION OF TOTAL VARIANCE, WHICH WILL BE VERY USEFUL IN SELECTING NUMBER OF PRINCIPAL COMPONENTS.

RESULT 4 THE TOTAL VARIANCE EQUALS TO THE SUM OF EIGENVALUES OF COVARIANCE MATRIX. THAT IS:

$$\sum_{i=1}^P \text{Var}(X_i) = \sum_{i=1}^P \lambda_i$$

Proof THE PROOF IS ALMOST TRIVIAL: cyclic property of trace.

$$\sum_{i=1}^P \text{Var}(X_i) = \text{tr}(\underline{\Sigma}) = \text{tr}(P \Lambda P') = \text{tr}(\Lambda P' P) = \text{tr}(\Lambda) = \sum_{i=1}^P \lambda_i$$