

# **Canonical Correlation Analysis**

**Recall:** The sample covariance matrix:

$$\mathbf{S}_{p \times p} = \begin{bmatrix} s_{11} & s_{12} & \cdots & s_{1p} \\ s_{12} & s_{11} & \cdots & s_{2p} \\ \vdots & \vdots & \ddots & \vdots \\ s_{1p} & s_{2p} & \cdots & s_{pp} \end{bmatrix}$$

where

$$s_{ik} = \frac{1}{n-1} \sum_{j=1}^n (x_{ij} - \bar{x}_i)(x_{kj} - \bar{x}_k)$$

The sample correlation matrix:

$$\mathbf{R}_{p \times p} = \begin{bmatrix} 1 & r_{12} & \cdots & r_{1p} \\ r_{12} & 1 & \cdots & r_{2p} \\ \vdots & \vdots & \ddots & \vdots \\ r_{1p} & r_{2p} & \cdots & 1 \end{bmatrix}$$

where

$$r_{ik} = \frac{s_{ik}}{\sqrt{s_{ii}} \sqrt{s_{kk}}} = \frac{\sum_{j=1}^n (x_{ij} - \bar{x}_i)(x_{kj} - \bar{x}_k)}{\sqrt{\sum_{j=1}^n (x_{ij} - \bar{x}_i)^2} \sqrt{\sum_{j=1}^n (x_{kj} - \bar{x}_k)^2}}$$

# Tests for independence and zero correlation

## Test for zero correlation (Independence between two normally distributed variables):

$$H_0: \rho_{ij} = 0 \text{ vs. } H_a: \rho_{ij} \neq 0$$

The test statistic

$$t = r_{ij} \sqrt{\frac{(n-2)}{1-r_{ij}^2}}$$

If independence is true then the test statistic  $t$  will have a  $t$ -distributions with  $\nu = n - 2$  degrees of freedom.

The test is to reject *independence* if:

$$|t| > t_{\alpha/2}^{(n-2)}$$

**Test for specific correlation** ( $H_0: \rho = \rho_0$  vs.  $H_a: \rho \neq \rho_0$ )

Large sample test statistic:

$$z = \frac{\frac{1}{2} \ln \left( \frac{1+r}{1-r} \right) - \frac{1}{2} \ln \left( \frac{1+\rho_0}{1-\rho_0} \right)}{\sqrt{\frac{1}{n-3}}}$$

If  $H_0$  is true the test statistic  $z$  will have approximately a standard normal distribution

We then *reject*  $H_0$  if:

$$|z| > z_{\alpha/2}$$

# **Partial Correlation**

Conditional Independence

## There are two groups of variables:

The first group  $\mathbf{X}^{(1)}$  has  $p$ -variate Normal distribution

The second group  $\mathbf{X}^{(2)}$  has  $q$ -variate Normal distribution

Let  $E(\mathbf{X}^{(1)}) = \boldsymbol{\mu}^{(1)}$  and  $E(\mathbf{X}^{(2)}) = \boldsymbol{\mu}^{(2)}$

$\text{cov}(\mathbf{X}^{(1)}) = \boldsymbol{\Sigma}_{11}$  and  $\text{cov}(\mathbf{X}^{(2)}) = \boldsymbol{\Sigma}_{22}$

$\text{cov}(\mathbf{X}^{(1)}, \mathbf{X}^{(2)}) = \boldsymbol{\Sigma}_{12}$

Then the two groups have a joint covariance matrix

$$\boldsymbol{\Sigma} = \begin{bmatrix} \boldsymbol{\Sigma}_{11} & \boldsymbol{\Sigma}_{12} \\ \boldsymbol{\Sigma}_{12}' & \boldsymbol{\Sigma}_{22} \end{bmatrix}$$



The matrix  $\Sigma_{2 \cdot 1} = \Sigma_{22} - \Sigma'_{12} \Sigma_{11}^{-1} \Sigma_{12}$

is called the *matrix of partial variances and covariances*.

The  $(i, j)^{\text{th}}$  element of the matrix  $\Sigma_{2 \cdot 1}$

$$\sigma_{ij \cdot 1, 2, \dots, p}$$

is called the *partial covariance* (variance if  $i = j$ ) between  $x_i$  and  $x_j$  given  $\mathbf{X}^{(1)} = x_1, \dots, x_p$ . Further,

$$\rho_{ij \cdot 1, 2, \dots, p} = \frac{\sigma_{ij \cdot 1, 2, \dots, p}}{\sqrt{\sigma_{ii \cdot 1, 2, \dots, p} \sigma_{jj \cdot 1, 2, \dots, p}}}$$

is called the *partial correlation* between  $x_i$  and  $x_j$  given  $\mathbf{X}^{(1)} = x_1, \dots, x_p$

## Sample partial covariance

Let:

$$S = \begin{bmatrix} S_{11} & S_{12} \\ S'_{12} & S_{22} \end{bmatrix}$$

denote the sample Covariance matrix

Let  $S_{2 \cdot 1} = S_{22} - S'_{12} S_{11}^{-1} S_{12}$

The  $(i, j)^{\text{th}}$  element of the matrix  $S_{2 \cdot 1}$

$$s_{ij \cdot 1, 2, \dots, p}$$

is called the *sample partial covariance* (variance if  $i = j$ )  
between  $x_i$  and  $x_j$  given  $x_1, \dots, x_p$

Also,

$$r_{ij \cdot 1, 2, \dots, p} = \frac{s_{ij \cdot 1, 2, \dots, p}}{\sqrt{s_{ii \cdot 1, 2, \dots, p} s_{jj \cdot 1, 2, \dots, p}}}$$

is called the *sample partial correlation* between  $x_i$  and  $x_j$   
given  $x_1, \dots, x_p$

## Test for zero partial correlation

(Conditional independence between a two variables given a set of  $p$  variables)

The test statistic

$$t = r_{ij \cdot x_1, \dots, x_p} \sqrt{\frac{(n - p - 2)}{1 - r_{ij \cdot x_1, \dots, x_p}^2}}$$

$r_{ij \cdot x_1, \dots, x_p}$  = the partial correlation between  $y_i$  and  $y_j$  given  $x_1, \dots, x_p$ .

If independence is true then the test statistic  $t$  will have a  $t$  - distributions with  $\nu = n - p - 2$  degrees of freedom.

The test is to reject *independence* if:

$$|t| > t_{\alpha/2}^{(n-p-2)}$$

## Large Sample Test for a specific partial correlation:

$$H_0 : \rho_{ij.x_1, \dots, x_p}^0 = \rho_{ij.x_1, \dots, x_p}^0$$

The test statistic

$$z = \frac{\frac{1}{2} \ln \left( \frac{1 + r_{ij.x_1, \dots, x_p}^0}{1 - r_{ij.x_1, \dots, x_p}^0} \right) - \frac{1}{2} \ln \left( \frac{1 + \rho_{ij.x_1, \dots, x_p}^0}{1 - \rho_{ij.x_1, \dots, x_p}^0} \right)}{\sqrt{\frac{1}{n - p - 3}}}$$

If  $H_0$  is true the test statistic  $z$  will have approximately a Standard Normal distribution

We then ***reject***  $H_0$  if:  $|z| > z_{\alpha/2}$

# **The Multiple Correlation Coefficient**

Testing independence between a single variable and a group of variables

## Definition:

Suppose  $\begin{bmatrix} X \\ y \end{bmatrix}$  has  $(p + 1)$ -variate Normal distribution

with covariance matrix

$$\Sigma = \begin{bmatrix} \Sigma_X & \Sigma_{Xy} \\ \Sigma_{Xy} & \sigma_y \end{bmatrix}$$

We are interested if the variable  $y$  is independent of the vector  $X$

The **multiple correlation coefficient** is the maximum correlation between  $y$  and a linear combination of the components of  $X$

It can be shown that the **multiple correlation coefficient** is:

$$\rho_{y|x_1, \dots, x_p} = \sqrt{\frac{\Sigma_{Xy} \Sigma_X^{-1} \Sigma'_{Xy}}{\sigma_y}}$$

**The sample multiple correlation coefficient:**

Let the sample correlation matrix be

$$S = \begin{bmatrix} S_X & S_{Xy} \\ S_{Xy} & s_y \end{bmatrix}$$

Then the sample multiple correlation coefficient is

$$r_{y|x_1, \dots, x_p} = \sqrt{\frac{S_{Xy} S_X^{-1} S'_{Xy}}{s_y}}$$



We are interested if the variable  $y$  is independent of the vector  $X$

That is,

$$H_0: \rho_{y|x_1, \dots, x_p} = \sqrt{\frac{\Sigma_{Xy} \Sigma_X^{-1} \Sigma'_{Xy}}{\sigma_y}} = 0$$

The test statistic is:

$$F = \frac{n - p - 1}{p} \frac{r_{y|x_1, \dots, x_p}^2}{1 - r_{y|x_1, \dots, x_p}^2} = \frac{n - p - 1}{p} \frac{S_{Xy} S_X^{-1} S'_{Xy}}{s_y - S_{Xy} S_X^{-1} S'_{Xy}}$$

If independence is true then the test statistic  $F$  will have an  $F$ -distributions with  $\nu_1 = p$  degrees of freedom in the numerator and  $\nu_2 = n - p - 1$  degrees of freedom in the denominator

The test is to reject *independence* if:  $F > F_\alpha(p, n - p - 1)$

# **Canonical Correlation Analysis**

# The problem

Quite often when one has collected data on several variables.

The variables are grouped into two (or more) sets of variables and the researcher is interested in whether one set of variables is independent of the other set.

In addition if it is found that the two sets of variates are dependent, it is then important to describe and understand the nature of this dependence.

The appropriate statistical procedure in this case is called *Canonical Correlation Analysis*.

# Canonical Correlation: An Example

In the following study the researcher was interested in whether specific instructions on how to *relax* when taking tests and how to increase *Motivation* , would affect performance on standardized achievement tests

- Reading,
- Language and
- Mathematics

A group of 65 third- and fourth-grade students were rated after the instruction and immediately prior taking the Scholastic Achievement tests on:

- how **relaxed** they were ( $X_1$ ) and
- how **motivated** they were ( $X_2$ ).

In addition data was collected on the three achievement tests

- Reading ( $Y_1$ ),
- Language ( $Y_2$ ) and
- Mathematics ( $Y_3$ ).

The data were tabulated on the next page

	Relaxation	Motivation	Reading	Language	Math		Relaxation	Motivation	Reading	Language	Math
Case	$X_1$	$X_2$	$Y_1$	$Y_2$	$Y_3$	Case	$X_1$	$X_2$	$Y_1$	$Y_2$	$Y_3$
1	7	14	311	436	154	34	40	20	362	416	107
2	43	25	501	455	765	35	40	18	596	592	622
3	32	21	507	473	702	36	35	17	431	346	493
4	17	12	453	392	401	37	33	17	361	414	404
5	23	12	419	337	284	38	40	27	663	451	651
6	10	16	545	538	414	39	31	15	569	462	398
7	22	21	509	512	491	40	29	19	699	622	478
8	13	19	320	308	517	41	37	16	187	223	221
9	31	21	357	296	496	42	21	23	1132	839	1044
10	24	26	485	372	685	43	24	15	457	410	400
11	26	21	811	748	902	44	19	14	413	448	520
12	35	20	367	436	393	45	33	22	569	605	615
13	24	17	242	349	137	46	19	19	650	685	440
14	20	8	237	140	331	47	26	22	424	427	482
15	38	27	417	648	618	48	20	15	475	604	742
16	32	19	429	446	458	49	22	21	519	612	446
17	14	11	555	579	438	50	37	22	338	463	327
18	24	12	599	497	414	51	41	28	674	613	534
19	38	25	403	383	606	52	29	35	381	624	565
20	30	8	550	324	674	53	25	12	199	171	316
21	22	25	377	496	242	54	27	21	577	523	699
22	36	28	671	585	710	55	22	20	425	466	402
23	3	22	498	488	481	56	4	11	392	192	354
24	44	28	477	583	260	57	27	22	401	520	558
25	24	25	609	413	670	58	28	23	321	410	460
26	33	18	521	522	716	59	33	20	682	433	743
27	24	21	495	645	491	60	33	24	719	727	1052
28	28	20	400	555	624	61	31	33	672	705	650
29	34	7	258	175	276	62	20	11	366	309	537
30	39	20	466	541	348	63	26	25	581	558	386
31	7	19	709	757	589	64	23	10	681	530	581
32	13	17	586	472	492	65	30	22	1019	917	880
33	32	18	418	361	428						

# **Definition:** (Canonical variates and Canonical correlations)

Let  $\mathbf{X}^{(1)}$  have  $p$ -variate Normal distribution  
and  $\mathbf{X}^{(2)}$  have  $q$ -variate Normal distribution

Let  $E(\mathbf{X}^{(1)}) = \boldsymbol{\mu}^{(1)}$  and  $E(\mathbf{X}^{(2)}) = \boldsymbol{\mu}^{(2)}$

$\text{cov}(\mathbf{X}^{(1)}) = \boldsymbol{\Sigma}_{11}$  and  $\text{cov}(\mathbf{X}^{(2)}) = \boldsymbol{\Sigma}_{22}$

$\text{cov}(\mathbf{X}^{(1)}, \mathbf{X}^{(2)}) = \boldsymbol{\Sigma}_{12}$

Then the two groups have a joint covariance matrix

$$\boldsymbol{\Sigma} = \begin{bmatrix} \boldsymbol{\Sigma}_{11} & \boldsymbol{\Sigma}_{12} \\ \boldsymbol{\Sigma}_{12}' & \boldsymbol{\Sigma}_{22} \end{bmatrix}$$

**derivation:** ( 1<sup>st</sup> pair of Canonical variates and Canonical correlation)

Let  $U = \mathbf{a}'\mathbf{X}^{(1)}$

and  $V = \mathbf{b}'\mathbf{X}^{(2)}$

We want to find  $U$  and  $V$  such that they achieve the maximum correlation  $\phi_1 = \text{corr}(U, V)$

Then the solution  $U_1$  and  $V_1$  are called the first pair of *canonical variates* and  $\phi_1$  is called the first *canonical correlation coefficient*.



**derivation:** ( 1<sup>st</sup> pair of Canonical variates and Canonical correlation)

We have that:

$$\text{Var}(U) = \mathbf{a}'\text{cov}(\mathbf{X}^{(1)})\mathbf{a} = \mathbf{a}'\mathbf{\Sigma}_{11}\mathbf{a}$$

$$\text{Var}(V) = \mathbf{b}'\text{cov}(\mathbf{X}^{(2)})\mathbf{b} = \mathbf{b}'\mathbf{\Sigma}_{22}\mathbf{b}$$

$$\text{Cov}(U, V) = \mathbf{a}'\text{cov}(\mathbf{X}^{(1)}, \mathbf{X}^{(2)})\mathbf{b} = \mathbf{a}'\mathbf{\Sigma}_{12}\mathbf{b}$$

Hence:

$$\text{Corr}(U, V) = \frac{\mathbf{a}'\mathbf{\Sigma}_{12}\mathbf{b}}{\sqrt{\mathbf{a}'\mathbf{\Sigma}_{11}\mathbf{a}}\sqrt{\mathbf{b}'\mathbf{\Sigma}_{22}\mathbf{b}}}$$

Thus we want to choose  $b$  and  $\mathbf{b}$

so that

$$\rho_{UV} = \frac{\mathbf{a}'\Sigma_{12}\mathbf{b}}{\sqrt{\mathbf{a}'\Sigma_{11}\mathbf{a}}\sqrt{\mathbf{b}'\Sigma_{22}\mathbf{b}}} \text{ is at a maximum}$$

### **Solution:**

The first pair of *canonical variates*

$$U_1 = \mathbf{a}_1'\mathbf{X}^{(1)} \text{ and } V_1 = \mathbf{b}_1'\mathbf{X}^{(2)}$$

are found by finding  $\vec{a}_1$  and  $\vec{b}_1$ , eigenvectors of the matrices

$$\Sigma_{11}^{-1/2}\Sigma_{12}\Sigma_{22}^{-1}\Sigma_{21}\Sigma_{11}^{-1/2} \text{ and } \Sigma_{22}^{-1/2}\Sigma_{21}\Sigma_{11}^{-1}\Sigma_{12}\Sigma_{22}^{-1/2}$$

associated with the largest eigenvalue (same for both matrices)

The largest eigenvalue of the two matrices is the square of the first canonical correlation coefficient  $\phi_1$

$$\phi_1 = \sqrt{\text{the largest eigenvalue of } \Sigma'_{12}\Sigma_{11}^{-1}\Sigma_{12}\Sigma_{22}^{-1}}$$

## The remaining canonical variates and canonical correlation coefficients

The second pair of *canonical variates*

$$U_2 = \mathbf{a}_2' \mathbf{X}^{(1)}$$

and  $V_2 = \mathbf{b}_2' \mathbf{X}^{(2)}$

are found by finding  $\vec{a}_2$  and  $\vec{b}_2$ , so that

1.  $(U_2, V_2)$  are independent of  $(U_1, V_1)$ .
2. The correlation between  $U_2$  and  $V_2$  is maximized

The correlation,  $\phi_2$ , between  $U_2$  and  $V_2$  is called the *second canonical correlation coefficient*.

The  $i^{th}$  pair of *canonical variates*

$$U_i = \mathbf{a}_i' \mathbf{X}^{(1)}$$

and  $V_i = \mathbf{b}_i' \mathbf{X}^{(2)}$

are found by finding  $\vec{a}_i$  and  $\vec{b}_i$ , so that

1.  $(U_i, V_i)$  are independent of  $(U_1, V_1), \dots, (U_{i-1}, V_{i-1})$ .
2. The correlation between  $U_i$  and  $V_i$  is maximized

The correlation,  $\phi_i$ , between  $U_i$  and  $V_i$  is called the  $i^{th}$  *canonical correlation coefficient*.

It can be shown that  $\mathbf{a}_2$  and  $\mathbf{b}_2$   
are eigenvectors of the matrices

$$\Sigma_{11}^{-1/2} \Sigma_{12} \Sigma_{22}^{-1} \Sigma_{21} \Sigma_{11}^{-1/2} \text{ and } \Sigma_{22}^{-1/2} \Sigma_{21} \Sigma_{11}^{-1} \Sigma_{12} \Sigma_{22}^{-1/2}$$

associated with the  $2^{nd}$  largest eigenvalue (same for both matrices)

The  $2^{nd}$  largest eigenvalue of the two matrices is the square of the  $2^{nd}$  canonical correlation coefficient  $\phi_2$

$$\begin{aligned} \phi_2 &= \sqrt{\text{the } 2^{nd} \text{ largest eigenvalue of } \Sigma'_{12} \Sigma_{11}^{-1} \Sigma_{12} \Sigma_{22}^{-1}} \\ &= \sqrt{\text{the } 2^{nd} \text{ largest eigenvalue of } \Sigma_{22}^{-1} \Sigma'_{12} \Sigma_{11}^{-1} \Sigma_{12}} \end{aligned}$$

## continuing

Coefficients for the  $i^{th}$  pair of canonical variates,  $\vec{a}_i$  and  $\vec{b}_i$  are eigenvectors of the matrices

$$\Sigma_{11}^{-1/2} \Sigma_{12} \Sigma_{22}^{-1} \Sigma_{21} \Sigma_{11}^{-1/2} \text{ and } \Sigma_{22}^{-1/2} \Sigma_{21} \Sigma_{11}^{-1} \Sigma_{12} \Sigma_{22}^{-1/2}$$

associated with the  $i^{th}$  largest eigenvalue (same for both matrices)

The  $i^{th}$  largest eigenvalue of the two matrices is the square of the  $i^{th}$  canonical correlation coefficient  $\phi_i$

$$\begin{aligned} \phi_i &= \sqrt{\text{the } i^{th} \text{ largest eigenvalue of } \Sigma'_{12} \Sigma_{11}^{-1} \Sigma_{12} \Sigma_{22}^{-1}} \\ &= \sqrt{\text{the } i^{th} \text{ largest eigenvalue of } \Sigma_{22}^{-1} \Sigma'_{12} \Sigma_{11}^{-1} \Sigma_{12}} \end{aligned}$$

# Example

## Variables

- **relaxation Score** ( $X_1$ )
- **motivation score** ( $X_2$ ).
- **Reading** ( $Y_1$ ),
- **Language** ( $Y_2$ ) and
- **Mathematics** ( $Y_3$ ).

# Summary Statistics

## UNIVARIATE SUMMARY STATISTICS

-----

VARIABLE	MEAN	STANDARD DEVIATION
1 Relax	26.87692	9.50412
2 Mot	19.41538	5.83066
3 Read	499.03077	172.25508
4 Lang	485.83077	156.08957
5 Math	512.52308	195.18614

## CORRELATIONS

-----

		Relax	Mot	Read	Lang	Math
		1	2	3	4	5
Relax	1	1.000				
Mot	2	0.391	1.000			
Read	3	0.002	0.280	1.000		
Lang	4	0.050	0.510	0.781	1.000	
Math	5	0.127	0.340	0.713	0.556	1.000



# Canonical Correlation Statistics

EIGENVALUE	CANONICAL CORRELATION	NUMBER OF EIGENVALUES	BARTLETT'S TEST FOR REMAINING EIGENVALUES		
			CHI- SQUARE	D.F.	TAIL PROB.
			27.86	6	0.0001
0.35029	0.59186	1	1.56	2	0.4586
0.02523	0.15885				

BARTLETT'S TEST ABOVE INDICATES THE NUMBER OF CANONICAL VARIABLES NECESSARY TO EXPRESS THE DEPENDENCY BETWEEN THE TWO SETS OF VARIABLES. THE NECESSARY NUMBER OF CANONICAL VARIABLES IS THE SMALLEST NUMBER OF EIGENVALUES SUCH THAT THE TEST OF THE REMAINING EIGENVALUES IS NON-SIGNIFICANT. FOR EXAMPLE, IF A TEST AT THE .01 LEVEL WERE DESIRED, THEN 1 VARIABLES WOULD BE CONSIDERED NECESSARY. HOWEVER, THE NUMBER OF CANONICAL VARIABLES OF PRACTICAL VALUE IS LIKELY TO BE SMALLER.