

MAT 4375

Multivariate Statistical Methods

Study Guide



Fall 2024

Vectors and Matrices

- A matrix is a rectangular array of scalar elements

$$A_{p \times q} = (a_{ij}) = \begin{pmatrix} a_{11} & a_{12} & \dots & a_{1q} \\ a_{21} & a_{22} & \dots & a_{2q} \\ \vdots & \vdots & \ddots & \vdots \\ a_{p1} & a_{p2} & \dots & a_{pq} \end{pmatrix}$$

$p \times q$ dimension of matrix

P: number of rows

Q: number of columns

Row vector: $a_{1j} = (a_{11}, a_{12}, a_{13}, \dots, a_{1q})$

Column vector:

$$a_{1i} = \begin{pmatrix} a_{11} \\ a_{21} \\ a_{31} \\ \vdots \\ a_{p1} \end{pmatrix}$$

Elementary Matrix Operations

Sum and difference: $A+B$ or $A-B$, done by doing element-wise addition and subtraction

Multiplication (written AB)

$$AB = \left(\sum_k a_{ik} b_{kj} \right)$$

- For multiplication to work, we need the matrices to be conformable

Inner product: Product of row vector and column vector, leads to scalar

Outer product: product of column vector and row vector, leads to matrix

Transpose: interchanging rows and columns of A , denoted A'

A square matrix has same rows and columns. A square matrix A is symmetric if $A=A'$

Properties of transpose:

- $(A')' = A$
- $(A+B)' = A'+B'$
- $(AB)' = B'A'$
- $(cA)' = cA'$

A diagonal matrix is a square matrix with zero off-diagonal elements

Upper triangular matrices have zero elements below the diagonal and lower triangular matrices have zero elements above the diagonal

Kronecker Product:

$$A = (a_{ij}) : p \times q$$

$$B = (b_{ij}) : m \times n$$

$$A \otimes B = \begin{pmatrix} a_{11}B & a_{12}B & \dots & a_{1q}B \\ a_{21}B & a_{22}B & \dots & a_{2q}B \\ \vdots & \vdots & \ddots & \vdots \\ a_{p1}B & a_{p2}B & \dots & a_{pq}B \end{pmatrix}, Pm \times Qn$$

Properties of Kronecker Product:

$$(A+B) \otimes C = (A \otimes C) + (B \otimes C)$$

$$A \otimes (B+C) = (A \otimes B) + (A \otimes C)$$

$$- A \otimes B \neq B \otimes A$$

$$(A \otimes B)(C \otimes D) = AC \otimes BD$$

$$-(A \otimes B)' = A' \otimes B'$$

Vectorization, denoted by $\text{vec}(A)$, is done by stacking the columns of A and we form a column vector

$$A = \begin{pmatrix} a & b \\ c & d \end{pmatrix}$$

$$\text{vec}(A) = \begin{pmatrix} a \\ c \\ b \\ d \end{pmatrix}$$

Trace of Matrix

- Sum of diagonals of matrix
 $\text{trace}(A) = \sum_i a_{ii}$
- Properties of trace
 - $\text{trace}(A') = \text{trace}(A)$
 - $\text{trace}(A+B) = \text{trace}(A) + \text{trace}(B)$
 - $\text{trace}(AB) = \text{trace}(BA)$
 - $\text{trace}(AA') = \text{trace}(A'A)$
 - $\text{trace}(ABC) = \text{trace}(CAB) = \text{trace}(BCA)$
 - $\text{trace}(A \otimes B) = \text{trace}(A)\text{trace}(B)$

Determinant of Matrix

- For 1x1 matrix, or scalar, the determinant is the value itself
- For 2x2 matrix $A = \begin{pmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{pmatrix}$, the determinant, denoted by $|A|$ is given by

$$|A| = \begin{vmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{vmatrix} = a_{11}a_{22} - a_{12}a_{21}$$

For 3x3 matrix

$$A = \begin{pmatrix} a_{11} & a_{12} & a_{13} \\ a_{21} & a_{22} & a_{23} \\ a_{31} & a_{32} & a_{33} \end{pmatrix}$$

The determinant is calculated as

$$a_{11} \begin{vmatrix} a_{22} & a_{23} \\ a_{32} & a_{33} \end{vmatrix} - a_{12} \begin{vmatrix} a_{21} & a_{23} \\ a_{31} & a_{33} \end{vmatrix} + a_{13} \begin{vmatrix} a_{21} & a_{22} \\ a_{31} & a_{32} \end{vmatrix}$$

General formula

$$\det(A) = \sum_{i=1}^n (-1)^{i+j} a_{ij} |A_{ij}| \text{ for all } j$$

or

$$\det(A) = \sum_{j=1}^n (-1)^{j+i} a_{ij} |A_{ij}| \text{ for all } i$$

Properties

- $-|A'| = |A|$
- $-|AB| = |A||B| = |BA|$
- $-k|A| = k^p |A|, A \text{ p} \times p$
- $-|AA'| \geq 0$

Positive semi-definite matrix: a symmetric matrix A is called positive semi-definite if $\bar{x}'A\bar{x} \geq 0$ for all $\bar{x} \in \mathbb{R}^p$ and $\bar{x}'A\bar{x} = 0$ for some nonzero \bar{x} .

A positive semi-definite matrix can be written as XX' for some matrix X , hence $|A| \geq 0$. Positive definite matrix: a symmetric matrix A is called positive definite if $\bar{x}'A\bar{x} > 0$ for any non-zero $\bar{x} \in \mathbb{R}^p$ vector. Again, A can be written as XX' , but in this case X is nonsingular, $|A| > 0$

Rank of Matrix

Denoted $P(A)$, is the number of linearly independent rows of the matrix

Can use elementary row operations to get the row-echelon form (or reduced row-echelon form). The rank of the matrix is the number of non-zero rows of the reduced row-echelon form

A square matrix $A_{p \times p}$ is said to be of full rank if $|A|=p$. That is all rows are linearly independent

Full rank $\Leftrightarrow |A| > 0 \Leftrightarrow$ positive-definite

Inverse of Matrix

- A square matrix is invertible (or non-singular) if it is of full rank. Otherwise, it is singular
- For a non-singular matrix $A_{p \times p}$, there exists a unique matrix (inverse), often denoted by A^{-1} s.t. $AA^{-1} = I_p$
- Can obtain inverse of A using reduced row echelon form $(|AI|_p) \rightarrow (I_p | A^{-1})$
- Properties of inverse:
 - $(A')^{-1} = (A^{-1})'$, $(A^{-1})^{-1} = A$
 - $(AB)^{-1} = B^{-1}A^{-1}$, $(A \otimes B)^{-1} = A^{-1} \otimes B^{-1}$
 - $(kA)^{-1} = \frac{1}{k} A^{-1}$ for $k \neq 0$, $|A^{-1}| = \frac{1}{|A|}$

Generalized Inverses

- Denoted by A^- s.t. $AA^-A = A$
- Not usually unique
- Moore-Penrose inverse is unique:
 - Moore-Penrose Conditions
 - AA^- and A^-A are symmetric
 - $AA^-A = A$ and $A^-AA^- = A^-$

Null Matrix

- Denoted by A^0 s.t. $AA^0 = 0$
- $I - A(A'A)^{-1}A'$ is a null matrix for A

Eigenvalues and Eigenvectors

- They play important roles in multivariate analysis
- Eigenvalues of a square matrix $A_{p \times p}$ are roots of $|A - \lambda I_p| = 0$
- Eigenvector \bar{v} corresponding to eigenvalue λ of A can be found by solving $(A - \lambda I_p)\bar{v} = 0$
- Condition number: $\frac{\lambda_{\max}}{\lambda_{\min}}$
- Condition index: $\frac{\lambda_{\max}}{\lambda_{\min}}$
- A full rank matrix has all non-zero eigenvalues
- Large condition number means we have an ill-conditioned problem
- Determinant and Trace
 - $\text{trace}(A) = \sum_{i=1}^p \lambda_i$
 - $|A| = \prod_{i=1}^p \lambda_i$
- λ_i : eigenvalues of $A_{p \times p}$

Matrix Decompositions

- Singular Value Decomposition: For any matrix A , there exist two orthogonal matrices U and V and a diagonal matrix D such that $A = UDV'$.
 - Diagonal values of D are singular values of A
 - Columns of U and V are called left and right singular vectors
- For a symmetric matrix A , there exists an orthogonal matrix T and a diagonal matrix Λ such that $A = T\Lambda T'$. The diagonal elements of Λ are the eigenvalues of A and the columns of T are the corresponding eigenvectors. This is known as eigen or spectral decomposition
- Cholesky Decomposition: For a symmetric positive definite matrix A , we have $A = TT'$, T lower triangular

Quadratic forms

If we have a vector $\bar{x} \in \mathbb{R}^n$, then we have $\bar{x}'A\bar{x} \in \mathbb{R}$ is a quadratic form for some $A \in \mathbb{R}^{n \times n}$

If we want to write $\sum_{i=1}^n (x_i - \bar{x}_i)^2$ in matrix form, we use the following:

- \bar{x} is a vector of size n where all elements are equal to 1
- J_{nn} is an $n \times n$ matrix where all elements are equal to 1 (write as J for simplicity)
- $J^2 = nJ$
- If we let $J' = \frac{1}{n}J$, we get $J'^2 = J'$
- Define $C = I - \frac{1}{n}J = I - J'$
- We get $\bar{x}'C\bar{x} = \sum_{i=1}^n (x_i - \bar{x}_i)^2$

Some common Matrix Derivatives

If we have a column vector

$$\bar{a} = \begin{pmatrix} a_1 \\ a_2 \\ \vdots \\ a_n \end{pmatrix}$$

Then we have

$$\frac{\partial f(\bar{a})}{\partial \bar{a}} = \begin{pmatrix} \frac{\partial f(\bar{a})}{\partial a_1} \\ \frac{\partial f(\bar{a})}{\partial a_2} \\ \vdots \\ \frac{\partial f(\bar{a})}{\partial a_n} \end{pmatrix}$$

For $f(\bar{a}) = \bar{c}'\bar{a}$, where \bar{c} is a column vector of the same length as \bar{a} . We then get

$$\frac{\partial f(\bar{a})}{\partial \bar{a}} = \bar{c}$$

For $f(\bar{a}) = \bar{a}'B\bar{a}$, $B = (b_{ij}) = (b_{11} \ b_{12} \ \dots \ b_{1n})$ a symmetric matrix. We have

$$\frac{\partial f(\bar{a})}{\partial \bar{a}} = 2B'\bar{a} = 2Ba^*$$

Probability Distributions

- If X is a random variable, its cumulative distribution function (cdf) is $F_X(x) = P(X \leq x)$ for all $x \in (-\infty, \infty)$

Discrete Random Variables

A random variable X is called a discrete random variable if it assumes only finite or denumerably infinite number of values

$F_X(x)$ is a step function

The probability mass function of a discrete random variable X that takes values x_i , where $i=1, 2, 3, \dots$, is

$$P(X=x_i) = P(X=x_i) \text{ for all } i$$

Common Discrete Distributions: Discrete Uniform, Bernoulli, Binomial, Poisson, Geometric, Negative Binomial

Continuous Random Variables

A random variable is continuous if it takes on any real numbers within its support

A random variable is continuous if $F_X(x)$ is continuous

The probability density function $f_X(x)$ for a continuous random variable is a function such that

$$F_X(x) = \int_{-\infty}^x f_X(t)dt \text{ for } x \in (-\infty, \infty)$$

$$\text{or } \frac{d}{dx} F_X(x) = f_X(x)$$

$$\cdot f_X(x) \geq 0 \text{ for all } x$$

$$\int_{-\infty}^{\infty} f_X(x)dx = 1$$

Common Continuous Distributions: Normal, Chi-Square, Exponential, Gamma, Beta, Student t, F

Expectation and Moments

- For a continuous random variable X , the expected value is

$$E[X] = \int_{-\infty}^{\infty} x f_X(x) dx$$

- For a discrete random variable X , the expected value is

$$E[X] = \sum x_i P(x_i) \text{ for all } x_i$$

Also known as the mean of distribution

Is a measure of central tendency

Other measures of central tendency include median and mode

For a random variable X , the r th non-central moment is defined as $E[X^r]$

The expected value is therefore the first non-central moment

The r th central moment is defined as $E[(X-EX)]^r$

The variance of a random variable X is defined as its second central moment, or $E[(X-EX)^2]$

The square root of the variance is referred to as the standard deviation

Other measures of variability include interquartile range (IQR) and range

Normal Distribution

- Symmetric, unimodal

- Closed under linear combinations, conditioning, and marginalizing

- Most commonly used distributions are transformations of the normal distribution

- CLT

- Density function

$$f_X(x) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}} I_{(-\infty, \infty)}(x)$$

$-\infty < \mu < \infty, \sigma > 0$

- $E[X] = \mu$ and $\text{Var}[X] = \sigma^2$

Distribution characterized by its mean and variance, written $X \sim N(\mu, \sigma^2)$

- Density usually denoted $\phi_{\mu, \sigma^2}(x)$ and CDF $\Phi_{\mu, \sigma^2}(x)$

- $\frac{X-\mu}{\sigma} \sim N(0, 1)$ standard normal, pdf $\phi(x)$ and cdf $\Phi(x)$

Relationship of Normal with Other Distributions

- If $Z \sim N(0, 1)$, then $Z^2 \sim \chi^2(1)$

- If $Z_1, Z_2, \dots, Z_n \sim N(0, 1)$, then $\sum_{i=1}^n Z_i^2 \sim \chi^2(n)$

- If $X_1, X_2, \dots, X_n \sim N(\mu, \sigma^2)$, then we have $\frac{(n-1)S^2}{\sigma^2} \sim \chi^2(n-1)$

- If $Z \sim N(0, 1)$ and $V \sim \chi^2(n)$, then

$$\frac{Z}{\sqrt{V/n}} \sim t(n)$$

- If $X_1, X_2, \dots, X_n \sim N(\mu, \sigma^2)$, then

$$\frac{\bar{X}-\mu}{\frac{s}{\sqrt{n}}} \sim t(n-1)$$

- If $V \sim \chi^2(r)$ and $W \sim \chi^2(s)$, then

$$\frac{V}{r} \sim F(r, s)$$

Random Samples and Statistics

- A random sample is a sequence of independent and identically distributed random variables from a given distribution (i.e. $X_1, X_2, \dots, X_n \stackrel{iid}{\sim} f(x)$)

- A statistic is a function of a random sample that is observable

$$T = f(X_1, X_2, \dots, X_n)$$

$T = \bar{X} = \frac{1}{n} \sum_{i=1}^n X_i$ is a statistic

$T = \frac{1}{n} \sum_{i=1}^n (X_i - \mu)$ is not a statistic

- The sampling distribution of a statistic

i.e. $\bar{X} \sim N(\mu, \sigma^2/n)$

Estimation Methods

- Method of moments: equate sample and population moments of the same order, then solve for the estimators of parameters

- Maximum likelihood estimate: The likelihood function is

$$L(\theta) = \prod_{i=1}^n f(x_i; \theta)$$

Then we solve for $\hat{\theta}$ that maximizes $L(\theta)$

- Unbiased estimator: $E[\hat{\theta} - \theta] = 0$

- Confidence interval: We say that $[L, U]$ is a $100(1-\alpha)\%$ confidence interval for a parameter θ if $P(L < \theta < U) = 1 - \alpha$

Hypothesis Testing

- We usually formulate hypotheses

$$H_0: \theta = \theta_0$$

$$H_1: \theta \neq \theta_0 \text{ or } \theta > \theta_0 \text{ or } \theta < \theta_0$$

- Compute a test statistic T and the p-value is $P(T > t_c)$, where t_c is the critical value. Reject H_0 if the p-value is less than a significance level α

Errors

	Accept	Reject
H_0 True	✓	Type I Error
H_0 False	Type II Error	✓

Random Vectors

- Random vectors are defined as a vector of random variables $\vec{X} = (X_1, X_2, \dots, X_p)$ and are associated with joint distributions, which is multivariate in nature

- The vector of random variables might consist of different outcomes taken from the same individuals, or a single variable/outcome measured at multiple time points

- If $\vec{X} = (X_1, X_2, \dots, X_p)$ is a random vector, its joint distribution is given by

$$F_{\vec{X}}(x_1, x_2, \dots, x_p) = P(X_1 \leq x_1, X_2 \leq x_2, \dots, X_p \leq x_p)$$

for $-\infty < x_i < \infty$ and $i=1, 2, \dots, p$

- The expected value of a random vector is given by

$$E[\vec{X}] = (E[X_1], E[X_2], \dots, E[X_p])$$

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

- For a random vector $\vec{X} = (X_1, X_2, \dots, X_p)$ with expected value $\vec{\mu} = (\mu_1, \mu_2, \dots, \mu_p)$, the 2nd central moment, also known as the dispersion matrix, is given by

$$E[(\vec{X} - \vec{\mu})(\vec{X} - \vec{\mu})']$$

Note: if \vec{X} is a column vector, then it is given by

$$E[(\vec{X} - \vec{\mu})(\vec{X} - \vec{\mu})']$$

- Often denoted by

$$\Sigma_{\vec{X} \vec{X}'} = \begin{pmatrix} \sigma_{11} & \sigma_{12} & \dots & \sigma_{1p} \\ \sigma_{21} & \sigma_{22} & \dots & \sigma_{2p} \\ \vdots & \vdots & \ddots & \vdots \\ \sigma_{p1} & \sigma_{p2} & \dots & \sigma_{pp} \end{pmatrix}$$

where $\sigma_{ii} = E[(X_i - \mu_i)^2] = \sigma_i^2$ and $\sigma_{ij} = E[(X_i - \mu_i)(X_j - \mu_j)]$

- Correlation matrix

$$\rho = (\rho_{ij}) = \frac{E[(X_i - \mu_i)(X_j - \mu_j)]}{\sqrt{\sigma_i \sigma_j}}$$

- Both Σ and ρ are symmetric and positive semi-definite matrices, in most practical cases they're positive definite

- The covariance matrix is often estimated by the sample moments corresponding to each of its elements

- The unbiased estimate of Σ is referred to as the sample covariance matrix

Moment Generating Functions

- For a univariate random variable

$$M_x(t) = E[e^{tX}]$$

- The mgf uniquely identifies distributions

- We also have

$$\frac{d^r}{dt^r} M_x(t)|_{t=0} = E[X^r]$$

- The mgf of a random vector of length p (row vector) is defined as

$$M_{\vec{X}}(t) = M_{\vec{X}}(t_1, t_2, \dots, t_p) = E[e^{t_1 X_1 + t_2 X_2 + \dots + t_p X_p}]$$

- We can obtain univariate mgf's by

$$M_{X_i}(t_i) = M_{\vec{X}}(0, 0, \dots, t_i, 0, \dots, 0)$$

Random Matrices

- Matrix random variables are defined as

$$X_{\vec{p} \vec{q}} = \begin{pmatrix} X_{11} & X_{12} & \dots & X_{1q} \\ X_{21} & X_{22} & \dots & X_{2q} \\ \vdots & \vdots & \ddots & \vdots \\ X_{p1} & X_{p2} & \dots & X_{pq} \end{pmatrix}$$

where the elements X_{ij} 's are themselves random variables

- The expectation for matrix random variables is

$$E[\vec{X}] = \begin{pmatrix} E[X_{11}] & E[X_{12}] & \dots & E[X_{1q}] \\ E[X_{21}] & E[X_{22}] & \dots & E[X_{2q}] \\ \vdots & \vdots & \ddots & \vdots \\ E[X_{p1}] & E[X_{p2}] & \dots & E[X_{pq}] \end{pmatrix}$$

However, we run into problem when we try to consider higher order moments for matrix random variables

- The most common approach is to use the vectorized form of the matrix \vec{X} , where the covariance matrix is represented as $\Sigma \otimes \Psi$, where Σ represents the covariance for column vectors and Ψ represents the covariance for row vectors

Multivariate Bias

- Suppose $\hat{\theta} = (\hat{\theta}_1, \hat{\theta}_2, \dots, \hat{\theta}_p)$ is an estimator for a vector parameter $\theta = (\theta_1, \theta_2, \dots, \theta_p)$. Multivariate bias is defined as

$$E[\hat{\theta} - \theta] = (E[\hat{\theta}_1] - \theta_1, E[\hat{\theta}_2] - \theta_2, \dots, E[\hat{\theta}_p] - \theta_p)$$

- An estimator is referred to as an unbiased estimator if $E[\hat{\theta} - \theta] = \vec{0}$

- Although the definition of bias is quite natural, we face a challenge when we attempt to compare two vector estimators with respect to their respective bias vectors

- For each vector estimator, we can calculate

$$D = \sqrt{(E[\hat{\theta}] - \theta)^T (E[\hat{\theta}] - \theta)}$$

Compare the values for both estimators

- Can also use absolute bias across elements of the vector

Multivariate MSE

- For univariate estimators, it is defined as $E[(\hat{\theta} - \theta)^2]$ and can be rewritten as

$$MSE(\hat{\theta}) = E[(\hat{\theta} - E[\hat{\theta}])^2] + (E[\hat{\theta} - \theta])^2$$

$$= Var(\hat{\theta}) + Bias(\hat{\theta})^2$$

- For vector estimators, the MSE is defined as

$$MSE(\hat{\theta}) = E[(\hat{\theta} - \theta)'(\hat{\theta} - \theta)]$$

This is an outer product, hence $MSE(\hat{\theta})$ is a matrix.

- We can rewrite the MSE as

$$MSE(\hat{\theta}) = \Sigma_{\hat{\theta}} + Bias(\hat{\theta})' Bias(\hat{\theta})$$

$\Sigma_{\hat{\theta}}$: variance covariance matrix of $\hat{\theta}$

- If we want to compare 2 vector estimators, we vectorize the matrix and calculate the Euclidean distance to compare MSE for two vector estimators

- $MSE(\hat{\theta})$ involves covariances in addition to variances, hence difficult to interpret

- To overcome difficulty in interpretation, we can use another definition of MSE

$$MSE^*(\hat{\theta}) = (Var(\hat{\theta}_1) + Bias(\hat{\theta}_1)^2, \dots, Var(\hat{\theta}_p) + Bias(\hat{\theta}_p)^2)$$

This corresponds to the sum of the diagonal elements of $\Sigma_{\hat{\theta}}$ and the diagonal elements of $Bias(\hat{\theta})' Bias(\hat{\theta})$. Also equivalent to vector of element-wise MSEs

- Also easier to extend to matrix parameters

$$Var(\hat{\theta}_1) + Bias(\hat{\theta}_1)^2 \dots Var(\hat{\theta}_p) + Bias(\hat{\theta}_p)^2$$

$$Bias = \begin{pmatrix} E[\hat{\theta}_{11} - \theta_{11}] & \dots & E[\hat{\theta}_{1p} - \theta_{1p}] \\ \vdots & \ddots & \vdots \\ E[\hat{\theta}_{p1} - \theta_{p1}] & \dots & E[\hat{\theta}_{pp} - \theta_{pp}] \end{pmatrix}$$

- MSE for vector estimators are vectors. Use distance or average of element-wise MSEs to compare estimators

- Both bias and MSE for matrix estimators are matrices. We vectorize the matrices first and do the same as with vector estimators to compare estimators

Joint and Marginal Distributions

- If we have a random vector $\vec{X} = (X_1, X_2, \dots, X_p)$, the joint CDF is given by

$$F_{\vec{X}}(x_1, x_2, \dots, x_p) = P(X_1 \leq x_1, X_2 \leq x_2, \dots, X_p \leq x_p)$$

- The joint density function $f_{\vec{X}}(x_1, x_2, \dots, x_p)$ is defined in a way such that

$$F_{\vec{X}}(x_1, x_2, \dots, x_p) = \int_{-\infty}^{x_1} \int_{-\infty}^{x_2} \dots \int_{-\infty}^{x_p} f_{\vec{X}}(x_1, x_2, \dots, x_p) dx_1 dx_2 \dots dx_p$$

Similarly, the joint density function can be found from the CDF by differentiating w.r.t. all elements of the p-vector

- The marginal distribution of X_i in a p-variate random vector is obtained by taking p-1 integrals w.r.t. all the remaining random variables

Conditional Distributions and Independence

- Consider a p-variate vector $\vec{X} = (X_1, X_2, \dots, X_p)$, the conditional distribution of a sub-vector $\vec{X}' = (X_1, X_2, \dots, X_{i-1}, X_{i+1}, \dots, X_p)$ given X_i is

$$f_{\vec{X}'}|_{X_i}(x_1, \dots, x_{i-1}, x_{i+1}, \dots, x_p | x_i) = \frac{f_{\vec{X}}(x_1, x_2, \dots, x_p)}{f_{X_i}(x_i)}$$

Conditional distribution of a sub-vector of random variables conditioned on another sub-vector of random variables can also be defined in a similar way.

A vector of random variables $\vec{X} = (X_1, \dots, X_p)$ are independent if and only if the joint distribution factorizes. That is

$$F_{\vec{X}}(x_1, x_2, \dots, x_p) = \prod_{i=1}^p F_{X_i}(x_i)$$

The joint density also factorizes

$$f_{\vec{X}}(x_1, x_2, \dots, x_p) = \prod_{i=1}^p f_{X_i}(x_i)$$

Alternatively, random variables are independent iff the joint MGF factorizes

$$M_{\vec{X}}(\vec{t}) = \prod_{i=1}^p M_{X_i}(t_i)$$

We also have if a vector of random variables $\vec{X} = (Y_1, Y_2, \dots, Y_p)$ are independent, then

$$E[g(X_1)g(X_2)\dots g(X_p)] = \prod_{i=1}^p E[g(X_i)]$$

Independence also means that conditional distributions are the same as unconditional distributions

Properties of Expectations and Covariances

Consider a random vector $\vec{X} = (X_1, X_2, \dots, X_p)$ and let $E[\vec{X}] = \vec{\mu}$, then

$$E[A\vec{X}'] = AE[\vec{X}] = A\vec{\mu}'$$

for any non-random matrix $A: q \times p$, q any constant

Similarly

$$E[\vec{X}B] = E[\vec{X}]B = \vec{\mu}B$$

for any non-random matrix $B:p \times r$, r any constant

$$E[\vec{X}'] = E[\vec{X}]'$$

Let $Cov(\vec{X}) = \Sigma$, then

$$\text{Cov}(A\vec{X}') = ACov(\vec{X})A' = A\Sigma A'$$

for any non-random matrix $A: q \times p$

Similarly, we can show that

$$\text{Cov}(\vec{X}B) = B'\text{Cov}(\vec{X})B = B'\Sigma B$$

for any non-random matrix $B:p \times r$

$$\text{Cov}(\vec{X}') = \text{Cov}(\vec{X})$$

Covariance of a non-random vector is 0

Covariance matrix is always positive semi-definite

Multivariate Normal Distribution

Univariate normal distribution density

$$f_X(x) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{1}{2}(\frac{x-\mu}{\sigma})^2}$$

A random vector $\vec{X} = (X_1, X_2, \dots, X_p)$ has a multivariate normal distribution iff every linear combination of \vec{X} has a univariate normal distribution

Multivariate normal distribution can also be defined with the aid of the univariate standard normal distribution, and by taking advantage of the fact that the normal distribution belongs to the location-scale family

We first consider a random vector $\vec{U} = (U_1, U_2, \dots, U_p)$, $U_i \sim N(0, 1)$. We have

$$E[\vec{U}] = \vec{0} \text{ and } \text{Cov}(\vec{U}) = I_p$$

The distribution of \vec{U}' is the same as the joint distribution of the U_i 's and is given by

$$f_{\vec{U}}(u_1, u_2, \dots, u_p) = \prod_{i=1}^p f_{U_i}(u_i) = \left(\frac{1}{\sqrt{2\pi}}\right)^p e^{-\frac{1}{2}\sum_{i=1}^p u_i^2}$$

This can also be written as

$$\begin{aligned} & \left(\frac{1}{\sqrt{2\pi}}\right)^p e^{-\frac{1}{2}\text{trace}(\vec{u}\vec{u}')} \\ & = \left(\frac{1}{\sqrt{2\pi}}\right)^p e^{-\frac{1}{2}\text{trace}(\vec{u}'\vec{u})} \\ & = \left(\frac{1}{\sqrt{2\pi}}\right)^p e^{-\frac{1}{2}\text{trace}(\vec{u}'\vec{u})} \end{aligned}$$

This distribution is referred to as the standard p -variate normal distribution, often denoted by $\vec{U} \sim N_p(\vec{0}, I_p)$.

A p -dimensional random vector $\vec{X} = (X_1, X_2, \dots, X_p)$ is said to have a multivariate normal distribution if it has the same distribution as $\vec{U} + \vec{U}'B$, where $\vec{U} \sim N_p(\vec{0}, I_p)$, B is a non-singular and symmetric matrix of dimension $p \times p$ and $\Sigma = B'B = BB$.

The multivariate normal distribution is often denoted as $\vec{X} \sim N_p(\vec{\mu}, \Sigma)$ and its density function is given as

$$\left(\frac{1}{\sqrt{2\pi}}\right)^p |\Sigma|^{-\frac{1}{2}} e^{-\frac{1}{2}(\vec{X}-\vec{\mu})'\Sigma^{-1}(\vec{X}-\vec{\mu})'}$$

This can be derived using change-of-variables

B is referred to as the square root of Σ , denote it by $\Sigma^{\frac{1}{2}}$

The density function can also be written as

$$\left(\frac{1}{\sqrt{2\pi}}\right)^p |\Sigma|^{\frac{1}{2}} e^{-\frac{1}{2}\text{trace}(\Sigma(\vec{X}-\vec{\mu})')^2}$$

Using the trace property, we can also write

$$\left(\frac{1}{\sqrt{2\pi}}\right)^p |\Sigma|^{\frac{1}{2}} e^{-\frac{1}{2}\text{tr}(\Sigma(\vec{X}-\vec{\mu})')^2}$$

When $p=1$, the above reduces to univariate normal density

Properties of Multivariate Normal

Consider $\vec{X} \sim N_p(\vec{\mu}, \Sigma)$, we have

$$E[\vec{X}] = \vec{\mu}$$

$$\text{Cov}(\vec{X}) = \Sigma$$

We can write Σ as

$$\Sigma_{p,p} = \begin{pmatrix} \sigma_{11} & \sigma_{12} & \dots & \sigma_{1p} \\ \sigma_{21} & \sigma_{22} & \dots & \sigma_{2p} \\ \vdots & \vdots & \ddots & \vdots \\ \sigma_{p1} & \sigma_{p2} & \dots & \sigma_{pp} \end{pmatrix}$$

The multivariate normal distribution can also be defined using MGFs

MGF of univariate standard normal: $M_{U_i}(t) = e^{\frac{1}{2}t^2}$

If we have a random vector $\vec{U} = (U_1, U_2, \dots, U_p) \sim N_p(\vec{0}, I_p)$, its mgf is $M_{\vec{U}}(t) = e^{\frac{1}{2}\vec{t}'\vec{t}}$

$$\vec{t} = (t_1, t_2, \dots, t_n)$$

If we have $\vec{X} = \vec{\mu} + \vec{U}\Sigma$, we get

$$M_{\vec{X}}(t) = e^{\vec{\mu}'\vec{t} + \frac{1}{2}\vec{t}'\Sigma\vec{t}}$$

This is the mgf for $\vec{X} \sim N_p(\vec{\mu}, \Sigma)$

Normal distribution is closed under linear combination

Theorem: If $\vec{X} = (X_1, X_2, \dots, X_p) \sim N_p(\vec{\mu}, \Sigma)$ and $C:p \times m$, then

$$\vec{X}C \sim N_p(\vec{\mu}C, C'\Sigma C)$$

Can be used to prove normality of individual elements

If we have k independent p -variate normal random vectors $\vec{X}_i \sim N_p(\vec{\mu}_i, \Sigma_i)$ for $i=1, 2, \dots, k$ and $A_i:p \times p$. Then

$$\sum_{i=1}^k \vec{X}_i A_i \sim N_p\left(\sum_{i=1}^k \vec{\mu}_i A_i, \sum_{i=1}^k \vec{\Sigma}_i A_i A_i'\right)$$

Limit Theorems

X_1, X_2, \dots, X_n any random sample

$$E[X_i] = \mu \quad \text{Var}[X_i] = \sigma^2$$

Univariate Weak Law of Large Numbers:

$$\bar{X} \xrightarrow{P} \mu$$

Univariate Central Limit Theorem:

$$\frac{\sqrt{n}(\bar{X} - \mu)}{\sigma} \xrightarrow{d} N(0, 1)$$

Now, consider iid p -variate random vectors \vec{X}_i with mean $\vec{\mu}$ and covariance matrix Σ

Multivariate Weak Law of Large Numbers

$$\bar{X} = (X_1, X_2, \dots, X_p) \xrightarrow{P} \vec{\mu}$$

Multivariate CLT

$$\sqrt{n}(\bar{X} - \vec{\mu}) \xrightarrow{d} N_p(0, \Sigma)$$

Wishart Distribution

For univariate random variables, we have

$$X \sim N(0, 1) \Rightarrow X^2 \sim \chi^2(1)$$

$$X_1, \dots, X_n \stackrel{iid}{\sim} N(0, 1) \Rightarrow \sum_{i=1}^n X_i^2 \sim \chi^2(n)$$

$$X_1, \dots, X_n \stackrel{iid}{\sim} N(\mu, \sigma^2) \Rightarrow \frac{(n-1)S^2}{\sigma^2} \sim \chi^2(n-1)$$

The Wishart distribution is a multivariate extension of chi-squared distribution. Let

$$\vec{X} = \begin{pmatrix} X_1 \\ X_2 \\ \vdots \\ X_p \end{pmatrix} \sim N_p(\vec{\mu}, \Sigma)$$

Then we get that $W = \vec{X}'\vec{X}$ has the Wishart distribution, written as $w_p(I, \Sigma)$

If $\vec{X} \sim N_p(\vec{\mu}, \Sigma)$, then $W = \vec{X}'\vec{X} \sim w_p(I, \Sigma)$

If $\vec{X} \sim N_p(\vec{\mu}, \Sigma)$, then $W = \vec{X}'\vec{X}$ has the non-central Wishart distribution

If we let $X_1, X_2, \dots, X_n \stackrel{iid}{\sim} N_p(\vec{\mu}, \Sigma)$ be the columns of \vec{X} , we can write $\vec{X} \sim N_p(\vec{\mu}, \Sigma, I_n)$ where $\vec{\mu}' = (\mu_1, \mu_2, \dots, \mu_n)$

A random matrix $W:p \times n$ is said to have the Wishart distribution if and only if it can be written as $W = \vec{X}'\vec{X}$, where $\vec{X} \sim N_p(\vec{\mu}, \Sigma, I_n)$, Σ positive definite nonsingular matrix. This is written as $w_p(n, \Sigma, \Delta)$, $\Delta = \vec{\mu}'\vec{\mu}$ non-centrality parameter. When $\vec{\mu} = 0$, we have the central wishart distribution $w_p(n, \Sigma)$

We can write $W = XX' = \sum_{i=1}^n X_i X_i'$ $X_i \stackrel{iid}{\sim} N_p(\vec{\mu}, \Sigma)$

Elements of W are sums of independent random variables. Diagonal elements have $X_i'^2(n)$ distribution

Wishart density:

$$f_W(w) = c |w|^{(n-p)/2} |w|^{-\frac{1}{2}\text{trace}(\Sigma^{-1}w)}$$

$$c = \left(\frac{1}{2}\right)^{np} T_p\left(\frac{n}{2}\right)^{-1} T_p(\cdot) \text{ multivariate Gamma}$$

Properties:

If $W_1 \sim w_p(n, \Sigma, \Delta_1)$ and $W_2 \sim w_p(m, \Sigma, \Delta_2)$ and independent, then we have $W_1 + W_2 \sim w_p(n+m, \Sigma, \Delta_1 + \Delta_2)$

If $W \sim w_p(n, \Sigma, \Delta)$ and $A:r \times p$, then $AWA' \sim w_r(n, A\Sigma A', \Delta A A')$

Sampling from Multivariate Normal

If we let $X_1, X_2, \dots, X_n \stackrel{iid}{\sim} N_p(\vec{\mu}, \Sigma)$, we get $\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i \sim N_p(\vec{\mu}, \frac{1}{n}\Sigma)$

We can also get if $S = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})(X_i - \bar{X})'$, we get $(n-1)S \sim w_p(n-1, \Sigma)$

We also have that \bar{X} and S are jointly sufficient for $\vec{\mu}$ and Σ

Estimation of $\vec{\mu}$ and Σ

We can obtain that the method of moments estimators for $\vec{\mu}$ and Σ are $\hat{\vec{\mu}} = \bar{X}$ and $\hat{\Sigma} = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})(X_i - \bar{X})'$ by solving the system

$$E[X_i] = \bar{X} \text{ and } E[X_i X_i'] = \frac{1}{n} \sum_{i=1}^n X_i X_i'$$

If we want to do MLE, the likelihood is

$$L(\vec{\mu}, \Sigma | X) = \left(\frac{1}{2\pi}\right)^p |\Sigma|^{\frac{1}{2}} \exp\left(-\frac{1}{2}\sum_{i=1}^n (X_i - \bar{X})'\Sigma^{-1}(X_i - \bar{X})\right)$$

The log-likelihood is given by

$$\ell(\vec{\mu}, \Sigma | X) = -n\log(\sqrt{2\pi}) + \frac{n}{2} \log|\Sigma| - \frac{1}{2} \sum_{i=1}^n (X_i - \bar{X})'\Sigma^{-1}(X_i - \bar{X})$$

If we take derivatives wr.t. $\vec{\mu}$, we get

$$\frac{\partial \ell}{\partial \vec{\mu}} = -\frac{1}{2} \sum_{i=1}^n (X_i - \bar{X})$$

If we set this to 0, we get $\hat{\vec{\mu}} = \bar{X}$

To get the MLE for Σ , we need the following matrix derivatives

$$\frac{\partial |A|}{\partial A} = (A')^{-1}|A| \text{ and } \frac{\partial \text{trace}(AB)}{\partial B} = A'$$

Rewrite the log-likelihood as

$$\ell(\vec{\mu}, \Sigma | X) = -n\log(\sqrt{2\pi}) + \frac{n}{2} \log|\Sigma| - \frac{1}{2} \text{trace}((X - \bar{X})(X - \bar{X})')$$

We can get

$$\frac{\partial \ell}{\partial \Sigma^{-1}} = \frac{n}{2} \frac{1}{|I_{n-1}|} \cdot \Sigma \cdot |I_{n-1}| - \frac{1}{2} \sum_{i=1}^n (X_i - \bar{X})(X_i - \bar{X})'$$

If we set this to 0, we get

$$\hat{\Sigma} = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})(X_i - \bar{X})'$$

Note: $\hat{\Sigma}$ is not an unbiased estimator for Σ .

Hypothesis Testing for $\vec{\mu}$

For univariate case, if we have a sample from $N(\mu, \sigma^2)$, where σ^2 is fixed but unknown. If we want to test

$$H_0: \mu = \mu_0 \text{ vs. } H_1: \mu \neq \mu_0$$

We reject H_0 when

$$T = \left| \frac{\bar{X} - \mu_0}{\frac{1}{\sqrt{n}}} \right| > t_{\alpha/2}(n-1)$$

We also have $T^2 \sim F_{1, n-1}$, so we can also reject H_0 when $T^2 > F_{1, n-1}(1, n-1)$

For multivariate normal, if we want to test

$$H_0: \vec{\mu} = \vec{\mu}_0 \text{ vs. } H_1: \vec{\mu} \neq \vec{\mu}_0$$

we use the Hotelling's T^2 statistic

$$T^2 = n(\bar{X} - \vec{\mu}_0)' S^{-1}(\bar{X} - \vec{\mu}_0)$$

- The Hotelling's T^2 statistic is also the likelihood ratio test statistic
- The distribution for T^2 is not a standard distribution, so we need to scale it

Roy's Union-Intersection Test statistic

$$RF = \left(\frac{n-p}{p(n-p)} \right) T^2 \sim F(p, n-p)$$

We reject H_0 when $RF > F_{1-\alpha}(p, n-p)$

Two Sample Test

Let $X = (x_1, x_2, \dots, x_m)$, $Y = (y_1, y_2, \dots, y_n)$ random samples from two independent p -variate normal distributions with common covariance matrix Σ . Say $N_p(\mu_1, \Sigma)$ and $N_p(\mu_2, \Sigma)$, respectively

We want to test the following:

$$H_0: \mu_1 = \mu_2 \quad \text{vs.} \quad H_1: \mu_1 \neq \mu_2$$

For univariate case, we reject H_0 when

$$T = \sqrt{\frac{n_1 n_2}{n_1 + n_2}} \left| \frac{\bar{Y} - \bar{Y}}{S_p} \right| > t_{1-\alpha/2}(n_1 + n_2 - 2)$$

$$S_p^2 = \frac{(n_1-1)S_1^2 + (n_2-1)S_2^2}{n_1 + n_2 - 2}$$

S_1^2, S_2^2 sample variance for groups

We can also reject H_0 when

$$T^2 > F_{1-\alpha}(1, n_1 + n_2 - 2)$$

For multivariate case, we have the two sample Hotelling's T^2 statistic is

$$T^2 = \left(\frac{n_1 n_2}{n_1 + n_2} \right) (\bar{Y} - \bar{Y})' S_p^{-1} (\bar{Y} - \bar{Y})$$

$$S_p = \frac{(n_1-1)S_1 + (n_2-1)S_2}{n_1 + n_2 - 2}$$

S_1, S_2 sample covariance matrices for groups

Two-sample Roy's UIT

$$RF = \left(\frac{f-p+1}{fP} \right) T^2 \sim F(p, f-p+1)$$

$$f = n_1 + n_2 - 2$$

Reject H_0 when $RF > F_{1-\alpha}(p, f-p+1)$

Testing for Normality

Shapiro-Wilk's test: Consider a univariate random sample x_1, x_2, \dots, x_n and suppose y_1, y_2, \dots, y_n are the corresponding order statistics. The Shapiro-Wilk's test for normality is given by

$$W = \frac{\sum_{i=1}^n a_i y_i}{S}$$

S sample standard deviation
 a_i function of order statistic of standard normal distribution

Kolmogorov-Smirnov's test: uses the difference between the empirical distribution of two random samples, and reject when the difference is large

$$D = \sup_x |F_U(x) - F_W(x)|$$

reject when D is large

To test for normality, compare the empirical distribution corresponding to the data to standard normal

In practice, we usually use graphical methods, by themselves or combined with formal tests

Testing Multivariate Normality

- One way is to vectorize the random sample $X = (x_1, x_2, \dots, x_n)$.
- We have if $x_i \stackrel{iid}{\sim} N_p(\mu, \Sigma)$, then $X \sim N_{pn}(\mu\vec{1}, \Sigma, I_n)$. Then we have

$$\text{Vec}(X) \sim N_{pn}(\text{Vec}(\mu), I_p \otimes \Sigma)$$

- Then we can use univariate tests and traditional graphical methods like q-q plot to test for normality

- Limitation: Vectorized data does not necessarily consist of independent samples

- Another way is to use Small's graphical method. Let c_i be defined as

$$c_i = \frac{n}{(n-1)} (x_i - \bar{x})' S^{-1} (x_i - \bar{x})$$

- When $x_i \stackrel{iid}{\sim} N_p(\mu, \Sigma)$, it can be shown that c_i 's are approximately distributed as $\text{Beta}(\frac{n}{2}, \frac{n-1-p}{2})$.

We can use Kolmogorov-Smirnov's test as well as graphical approaches to test

ANOVA

- For comparing means between groups

$$\begin{array}{cccc} \text{Group 1} & \text{Group 2} & \cdots & \text{Group k} \\ y_{11} & y_{12} & \cdots & y_{1k} \\ y_{21} & y_{22} & \cdots & y_{2k} \\ \vdots & \vdots & \ddots & \vdots \\ y_{n_1,1} & y_{n_1,2} & \cdots & y_{n_1,k} \end{array}$$

- k groups with n_j individuals in j th group and $n = n_1 + n_2 + \dots + n_k$. y_{ij} is outcome for i th individual in j th group

- Can be represented as

$$y_{ij} = \mu_j + \varepsilon_{ij} \quad \varepsilon_{ij} \sim N(0, \sigma^2) \quad j=1, 2, \dots, k \quad i=1, 2, \dots, n$$

- To test hypotheses

$$H_0: \mu_1 = \mu_2 = \dots = \mu_k$$

$$H_1: \mu_k \neq \mu_j \text{ for at least one } k \neq j$$

We have

$$SST = \sum_{i,j} (y_{ij} - \bar{Y})^2$$

$$SSE = \sum_{i,j} (y_{ij} - \bar{y}_{\cdot j})^2$$

$$SSG = \sum_j n_j (\bar{y}_{\cdot j} - \bar{Y})^2$$

$$SST = SSE + SSG$$

The hypothesis is tested using

$$F = \frac{SSG/k-1}{SSE/n-k} \sim F_{k-1, n-k}$$

ANOVA Matrix Formulation

- Can be written as

$$Y_{1:n} = B_{1:k} X_{k:n} + E_{1:n}$$

$$Y = (y_{11}, y_{12}, \dots, y_{1k}, y_{21}, \dots, y_{nk})'$$

$$B = (B_{11}, B_{12}, \dots, B_{1k})'$$

$$X = \begin{pmatrix} 1 & 1 & \cdots & 1 \\ n_1 & n_2 & \cdots & n_k \\ \vdots & \vdots & \ddots & \vdots \\ 1 & 1 & \cdots & 1 \end{pmatrix}$$

$$E = (E_{11}, E_{12}, \dots, E_{1k}, E_{21}, \dots, E_{nk})'$$

- $\vec{1}_n$ and 0_n are row vectors of 1 s and 0 s with length n :

Matrix X referred to as design matrix

$$\text{We assume } E[\varepsilon] = 0, \text{Cov}(\varepsilon) = \sigma^2 I$$

Normality assumption: $\varepsilon \sim N_n(0, \sigma^2 I)$
This also means $\Sigma \sim N_n(B\Sigma, \sigma^2 I)$

- B has the same interpretation as in regression models in general
- When dummy variables are involved, we need to be careful with interpretation of elements in B

ANOVA Estimation

- We want to minimize

$$Q = (Y - BX)(Y - BX)'$$

- Then we get if we minimize Q , we get

$$\hat{B} = Y X' (X X')^{-1} \quad \hat{Y} = Y X' (X X')^{-1} X = Y H$$

- Residuals: $\hat{Y} - \hat{Y} = Y(I - H)$

MANOVA model

$$\begin{array}{cccc} \text{Group 1} & \text{Group 2} & \cdots & \text{Group k} \\ y_{11}, y_{12}, \dots, y_{1p} & y_{21}, y_{22}, \dots, y_{2p} & \cdots & y_{k1}, y_{k2}, \dots, y_{kp} \\ y_{11}, y_{12}, \dots, y_{1p} & y_{21}, y_{22}, \dots, y_{2p} & \cdots & y_{k1}, y_{k2}, \dots, y_{kp} \\ \vdots & \vdots & \ddots & \vdots \\ y_{n_1,1}, y_{n_1,2}, \dots, y_{n_1,p} & y_{n_2,1}, y_{n_2,2}, \dots, y_{n_2,p} & \cdots & y_{n_k,1}, y_{n_k,2}, \dots, y_{n_k,p} \end{array}$$

The model is given by: $Y = B X + E$

where

$$Y = \begin{pmatrix} y_{11} & y_{12} & \cdots & y_{1k} & \cdots & y_{1p} \\ y_{21} & y_{22} & \cdots & y_{2k} & \cdots & y_{2p} \\ \vdots & \vdots & \ddots & \vdots & \ddots & \vdots \\ y_{n_1,1} & y_{n_1,2} & \cdots & y_{n_1,k} & \cdots & y_{n_1,p} \\ y_{n_2,1} & y_{n_2,2} & \cdots & y_{n_2,k} & \cdots & y_{n_2,p} \\ \vdots & \vdots & \ddots & \vdots & \ddots & \vdots \\ y_{n_k,1} & y_{n_k,2} & \cdots & y_{n_k,k} & \cdots & y_{n_k,p} \end{pmatrix}'$$

$$B = \begin{pmatrix} p_1 & p_2 & \cdots & p_k \\ p_{11} & p_{12} & \cdots & p_{1k} \\ p_{21} & p_{22} & \cdots & p_{2k} \\ \vdots & \vdots & \ddots & \vdots \\ p_{n_1,1} & p_{n_1,2} & \cdots & p_{n_1,k} \end{pmatrix} \quad X = \begin{pmatrix} 1 & 0 & \cdots & 0 \\ 0 & 1 & \cdots & 0 \\ 0 & 0 & \cdots & 1 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & 1 \end{pmatrix}$$

$$E = \begin{pmatrix} \varepsilon_{11} & \varepsilon_{12} & \cdots & \varepsilon_{1k} & \cdots & \varepsilon_{1p} \\ \varepsilon_{21} & \varepsilon_{22} & \cdots & \varepsilon_{2k} & \cdots & \varepsilon_{2p} \\ \vdots & \vdots & \ddots & \vdots & \ddots & \vdots \\ \varepsilon_{n_1,1} & \varepsilon_{n_1,2} & \cdots & \varepsilon_{n_1,k} & \cdots & \varepsilon_{n_1,p} \\ \varepsilon_{n_2,1} & \varepsilon_{n_2,2} & \cdots & \varepsilon_{n_2,k} & \cdots & \varepsilon_{n_2,p} \\ \vdots & \vdots & \ddots & \vdots & \ddots & \vdots \\ \varepsilon_{n_k,1} & \varepsilon_{n_k,2} & \cdots & \varepsilon_{n_k,k} & \cdots & \varepsilon_{n_k,p} \end{pmatrix}'$$

- We assume the columns of E follow $N_p(0, \Sigma)$

- Can write as $E \sim N_{pn}(0, \Sigma, I_n)$

- We have $Y \sim N_{pn}(B\Sigma, \Sigma, I_n)$

MANOVA Inference

\hat{Y}	R
$C(X')$	$C(X')^2$

$$\hat{B} = Y X' (X X')^{-1} \quad \hat{Y} = \hat{B} X = Y X' (X X')^{-1} X = Y H$$

$$R = Y - \hat{Y} = Y(I - X'(X X')^{-1} X) = Y(I - H)$$

$X'(X X')^{-1} X$ symmetric and idempotent, hence a projection matrix

Projects Y on to the column space of X'

The matrix $I - X'(X X')^{-1} X$ projects onto the orthogonal complement of column space of X'

$\hat{\Sigma} = S_p = \frac{1}{n-k} R^T R$ is pooled sample variance-covariance matrix and is unbiased estimator of Σ

We also have $Y Y' = \hat{Y} \hat{Y}' + R R'$

\hat{Y} and R are independent

MANOVA Tests

Using vector formulation, the MANOVA model can be written as $y_{ij} = \mu_j + \varepsilon_{ij}$ where $\varepsilon_{ij} \sim N(0, \sigma^2)$

Matrix formulation: $Y = BX + E$, y_{ij} corresponds to the columns of Y , ε_{ij} corresponds to columns of E

Sample mean of group j : $\bar{y}_{\cdot j}$

Overall mean: \bar{Y}

Then we have that

$$\begin{aligned} SST &= \sum_{j=1}^k \sum_{i=1}^{n_j} (y_{ij} - \bar{y}_{\cdot j}) (y_{ij} - \bar{y}_{\cdot j})' \\ &= \sum_{j=1}^k \sum_{i=1}^{n_j} (y_{ij} - \bar{y}_{\cdot j}) (y_{ij} - \bar{y}_{\cdot j})' + \sum_{j=1}^k n_j (\bar{y}_{\cdot j} - \bar{Y}) (\bar{y}_{\cdot j} - \bar{Y})' \\ &= SSE + SSG \end{aligned}$$

We have

$$Y Y' = \hat{Y} \hat{Y}' + R R' \Rightarrow Y C Y' = \hat{Y} C \hat{Y}' + R C R' \quad \begin{matrix} SST \\ SSE \\ SSG \end{matrix}$$

We wish to test

$$H_0: \mu_1 = \mu_2 = \dots = \mu_k \quad \text{vs.} \quad H_1: \mu_j \neq \mu_k \text{ for atleast one } i \neq j$$

Under H_0 , MLE for Σ is

$$\hat{\Sigma} = \frac{1}{n} \sum_{j=1}^k \sum_{i=1}^{n_j} (y_{ij} - \bar{y}_{\cdot j}) (y_{ij} - \bar{y}_{\cdot j})'$$

$$\Rightarrow n \hat{\Sigma} = R C R'$$

Under H_0 UH, MLE for Σ is

$$\hat{\Sigma}_1 = \frac{1}{n} \sum_{j=1}^k \sum_{i=1}^{n_j} (y_{ij} - \bar{y}_{\cdot j}) (y_{ij} - \bar{y}_{\cdot j})'$$

$$\Rightarrow n \hat{\Sigma}_1 = R C R'$$

The Likelihood-Ratio test is given by

$$\frac{|R|}{|R_0|} = \frac{|SSE|}{|SST|} = \frac{|SSE|}{|SSG + SSE|}$$

$$H = SSG \quad E = SSE$$

$$\Lambda = \frac{|E|}{|E + H|} = \frac{1}{|I + H|} \quad (\text{Wilks Lambda})$$

Now, if we want to test $H_0: B=0$ vs $H_1: B \neq 0$

The LR test is then

$$\frac{|RR'|}{|YY'|} = \frac{|RR'|}{|RR + YY'|} = \frac{|SSE|}{|SSG + SSE|} = \frac{|SSE|}{|SST|}$$

Again, this is Wilks Lambda

$$\Lambda = \frac{1}{|I + H|}$$

This can also be written as

$$\Lambda = \frac{1}{\prod(1+\lambda_i)} \quad \lambda_i: \text{eigenvalues of } H E^{-1}$$

Other Tests for MANOVA

- Lawley-Hotelling trace test:

$$\text{Trace}(\mathbf{H}\mathbf{E}^{-1}) = \sum \lambda_i$$

reject for large values

- Pillai's Test : $\text{Trace}(\mathbf{H}(\mathbf{E}+\mathbf{H})^{-1}) = \sum \left(\frac{\lambda_i}{\lambda_i + \lambda_H} \right)$

- Roy's maximum root test: uses maximum eigenvalue of $\mathbf{H}\mathbf{E}^{-1}$ λ_{\max}

General Linear Hypothesis

- Test hypotheses:

$$H_0: G\mathbf{B}\mathbf{F} = 0 \quad \text{vs. } H_1: G\mathbf{B}\mathbf{F} \neq 0$$

G and F conformable matrices

- Wilks' Lambda:

$$\Lambda = \frac{|E|}{|E+H|} = \frac{1}{|I+\mathbf{H}\mathbf{E}^{-1}|}$$

$$H = G\hat{B}\hat{F}(F'(XX')^{-1}F)^{-1}F'\hat{B}'G'$$

$$E = GY(I - X'(XX')^{-1}X)Y'G = GRR'G'$$

Longitudinal Data

Group 1	Group 2	...	Group k
$y_{11}, y_{12}, \dots, y_{1p}$	$y_{21}, y_{22}, \dots, y_{2p}$...	$y_{k1}, y_{k2}, \dots, y_{kp}$
$y_{11}, y_{12}, \dots, y_{1p}$	$y_{21}, y_{22}, \dots, y_{2p}$...	$y_{k1}, y_{k2}, \dots, y_{kp}$
\vdots			
$y_{11}, y_{12}, \dots, y_{1p}$	$y_{21}, y_{22}, \dots, y_{2p}$...	$y_{k1}, y_{k2}, \dots, y_{kp}$

- Repeated measurements are from the same outcome at different time points
- There is often a functional relationship between the outcome and time
- MANOVA accounts for the correlation within individuals, but it assumes unstructured mean
- Time dependency often exists in most real life scenarios
- Repeated measurements can also be taken w.r.t. other continuous variables

GMANOVA Models

- GMANOVA arises as a linearly constrained MANOVA model
- Structure in the mean is induced in the model through a second design matrix
- GMANOVA accounts not only for the correlation between measurements across time, but it also allows us to model the response curves
- Often referred to as growth curve model (GCM), or multivariate bilinear regression model
- GMANOVA assumes the mean for each group to follow a polynomial function (of time) of degree $q-1$

$$E[Y|t] = \beta_{0,j} + \beta_{1,j}t + \dots + \beta_{q-1,j}t^{q-1}, \quad j=1,2,\dots,k$$

Note: Equal degree of polynomials assumed
Using matrix setup, the model can be formulated as

$$Y = ZBX + E$$

- Y : observation/outcome matrix

- B : parameter matrix

- Z : time dependency, within individual design matrix

- X : between individual design matrix

- $E \sim N_{p,n}(0, \Sigma, I_n)$ and $Y \sim N_{p,n}(ZBX, \Sigma, I_n)$

- It is assumed $p \geq q$ and $n \geq p(X) + p$

- $n = n$, that...+ n_k

We have the matrices are

$$Y = \begin{pmatrix} y_{11} & y_{12} & \dots & y_{1p} \\ y_{21} & y_{22} & \dots & y_{2p} \\ \vdots & \vdots & \ddots & \vdots \\ y_{k1} & y_{k2} & \dots & y_{kp} \end{pmatrix}$$

$$B = \begin{pmatrix} \beta_{0,1} & \beta_{1,1} & \dots & \beta_{q-1,1} \\ \beta_{0,2} & \beta_{1,2} & \dots & \beta_{q-1,2} \\ \vdots & \vdots & \ddots & \vdots \\ \beta_{0,k} & \beta_{1,k} & \dots & \beta_{q-1,k} \end{pmatrix}$$

$$Z = \begin{pmatrix} 1 & t_1 & t_1^2 & \dots & t_1^{q-1} \\ 1 & t_2 & t_2^2 & \dots & t_2^{q-1} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & t_p & t_p^2 & \dots & t_p^{q-1} \end{pmatrix}$$

$$X = \begin{pmatrix} 1 & n_1 & 0 & \dots & 0 \\ 0 & n_2 & 0 & \dots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \dots & 1 \end{pmatrix}$$

Inference for GMANOVA

- Consider GCM $Y = ZBX + E$
- First approach involves a linear transformation by pre-multiplying the model using the matrix $(Z'V^{-1}Z)^{-1}Z'V^{-1}$ where V is an arbitrary non-singular matrix

- The transformation leads to the MANOVA model

$$Y^* = BX + E^*$$

$$Y^* = (Z'V^{-1}Z)^{-1}Z'V^{-1}Y$$

$$E^* = (Z'V^{-1}Z)^{-1}Z'V^{-1}E$$

- Y^* and E^* normally distributed

- MLE of B is

$$\hat{B} = (Z'V^{-1}Z)^{-1}Z'V^{-1}YX'(XX')^{-1}$$

- Fitted values and residuals

$$\hat{Y} = Z\hat{B}X = Z(Z'V^{-1}Z)^{-1}Z'V^{-1}YX'(XX')^{-1}X$$

$$R = Y - \hat{Y} = Y - Z(Z'V^{-1}Z)^{-1}Z'V^{-1}YX'(XX')^{-1}X$$

- MLE of Z : $\hat{Z} = \frac{1}{n}RR'$
- When V is non-random, R is a linear combination of a multivariate normal random variable, and hence is normally distributed

Bilinear Projections

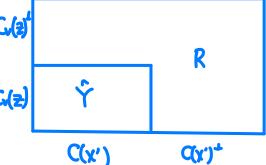
- If we look at the fitted values and residuals we obtained from Pothoff and Roy's transformations

$$Y = Z(Z'V^{-1}Z)^{-1}Z'V^{-1}YX'(XX')^{-1}X$$

$$= P_Z Y P_X$$

$$R = Y - P_Z Y P_X$$

$$P_Z = Z(Z'V^{-1}Z)^{-1}Z'V^{-1} P_X = X(XX')^{-1}X$$



- Decompositions of orthogonal complement to the design space has been done to facilitate usage and interpretation of residuals

- MLEs of parameters given by

$$\hat{B} = (Z'S^{-1}Z)^{-1}Z'S^{-1}YX'(XX')^{-1}$$

$$S = Y(1 - X'(XX')^{-1}X)Y$$

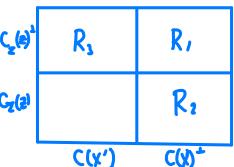
- Estimated mean trajectory

$$\hat{Y} = Z(Z'S^{-1}Z)^{-1}Z'S^{-1}YX'(XX')^{-1}Y$$

$$= P_{Z,S} Y P_X$$

$$P_{Z,S} = Z(Z'S^{-1}Z)^{-1}Z'S^{-1}$$

- The vector operator can be applied and the residual space can be decomposed into 3 orthogonal spaces



$$R_1 = (I - P_{Z,S})Y(1 - P_{Z,S})$$

$$R_2 = P_{Z,S}Y(1 - P_{Z,S})$$

$$R_3 = (I - P_{Z,S})Y P_X$$

Dimension Reduction

- One challenge with multivariate data and its analysis is difficulty in data visualization
- Understanding relationships between variables within the same data matrix or between data matrices is often challenging
- In some cases, even pairwise explorations are not feasible
- High-dimensionality can also cause a problem in subsequent confirmatory analyses
- Too many variables can cause curse of dimensionality
- Principal component analysis (PCA) is one of the statistical methods we use to overcome the challenges
- Exploratory in nature and is often used as dimension reduction and data visualization technique

Principal Component Analysis

- Involves projections/transformations where new sets of uncorrelated variables are created

- Main objective: reduce dimension of data while keeping much of the variation in the original data

- New variables are linear combinations of the original variables such that the variances of them are sequentially maximized

- Suppose $\mathbf{x} = (x_1, x_2, \dots, x_p)'$ is a random vector with $E[\mathbf{x}] = \mu$ and $Cov(\mathbf{x}) = \Sigma$, the new set of variables are

$$Z_1 = a_{11}x_1 + a_{12}x_2 + \dots + a_{1p}x_p = a_1' \mathbf{x}$$

$$Z_2 = a_{21}x_1 + a_{22}x_2 + \dots + a_{2p}x_p = a_2' \mathbf{x}$$

$$\vdots$$

$$Z_p = a_{p1}x_1 + a_{p2}x_2 + \dots + a_{pp}x_p = a_p' \mathbf{x}$$

- Z_i 's called principal components

- a_{ij} 's called loadings

- a_i 's called loading vectors

- Dimension reduction is achieved because most of the variation in the original data is often contained in the first few PCs, hence we often use the leading PCs in subsequent confirmatory analysis

- We want to transform the data such that the linear combinations maximize the variance sequentially.

- The first PC is obtained by maximizing its variance, while putting constraint $a_i'a_i = 1$

- Variance of first PC:

$$\text{Var}(z_i) = a_i' \Sigma a_i$$

- For the i th PC, we maximize the variance ($a_i' \Sigma a_i$) under the constraints $a_i'a_i = 1$ and $a_j'a_i = 0$ for any $j < i$

Obtaining Loading Vectors

- Σ can be decomposed as $T\Lambda T'$, where T is orthogonal and columns are eigenvectors of Σ , Λ is diagonal matrix with eigenvalues of Σ on the diagonals

- We have

$$\text{Var}(z_i) = \frac{a_i' \Sigma a_i}{a_i'a_i}$$

$$= \frac{a_i' \Lambda a_i}{a_i'a_i}$$

$$= \frac{b_i' \Lambda b_i}{b_i'b_i}$$

$$= \frac{\sum \lambda_i b_i^2}{\sum b_i^2} \leq \frac{\lambda_1 \sum b_i^2}{\sum b_i^2} = \lambda_1$$

λ_1 is maximum eigenvalue of Σ

- Equality is achieved when $a_i = b_1$, b_1 eigenvector of Σ corresponding to λ_1

- Variance of first PC is the largest eigenvalue and maximum variance is achieved by choosing the loading vector a_1 to be the normalizing eigenvector corresponding to λ_1

- Similarly, if we use $a_2 = b_2$, we get the variance to be λ_2 , second largest eigenvalue of Σ

- Then if $a_1 = b_1$, the variance is λ_1

- The loading vectors are orthogonal to each other, ensuring the PCs are uncorrelated

- We also have $\text{Var}(z_1) \geq \text{Var}(z_2) \geq \dots \geq \text{Var}(z_p)$

- The PCs preserve the total variation contained in the original data

$$\text{Var}(z_i) = \lambda_i; \quad \sum_{i=1}^p \text{Var}(z_i) = \text{trace}(\Sigma) = \sum_{i=1}^p \lambda_i$$

- In practice, we don't use all the PCs, hence there will be a loss of some information

- Proportion of total variance contained in k th PC is

$$PV_k = \frac{\lambda_k}{\sum_{i=1}^p \lambda_i}$$

- All of the above also valid for using correlation matrix. In practice, the correlation matrix is used more since variables often have different units of measurements and levels of variability

- Σ often unknown, so we use sample covariance matrix S or sample correlation matrix R

- PCs obtained using S or Σ are in most cases, different from R or P

- Total variation explained by PCs obtained from S or Σ are also different from those obtained using R or P

PCA of Random Sample

- Consider a random sample $\mathbf{x} = (x_1, x_2, \dots, x_p)$ where $E[\mathbf{x}_i] = \mu_i$ and $Cov(\mathbf{x}_i) = \Sigma_i$

- The new, transformed data is calculated from the PCs and is given by

$$Z = (z_1, z_2, \dots, z_p) = X\hat{P}$$

\hat{P} matrix of normalized eigenvectors of Σ

Cluster Analysis

- An exploratory statistical method used to uncover categories in the data, where individuals are homogenous within a group and are different between groups

- Often referred to as class discovery

- Unsupervised: group identifier often unknown and is not used in the analysis

- Assign individuals to groups (clusters) with the objective of maximizing homogeneity (similarity) within groups/clusters and maximizing difference between them

- Patterns in the data and potential clusters can sometimes be discovered through other exploratory analysis

- There may also be prior knowledge/evidence that motivates the study leading to class discovery

- Cluster analysis can also be simply used as a means of hypothesis generating

Distances

- The key in cluster analysis is similarity (homogeneity) and separation/heterogeneity (distance). We sometimes use the word dissimilarity to refer to distance matrices

- Most commonly used distance matrices are based on the Euclidean distance, and those based on the correlation or covariance matrix

- There are situations where one needs to use other similarity matrices

- Choice of an appropriate similarity/distance matrix that optimally quantifies similarities between individuals is, therefore, fundamental to performing cluster analysis

Some common Distances

- For any distance function d , we need
 - $d \geq 0$
 - $d(x_i, x_i) = 0$
- Let $x_j = (x_{j1}, x_{j2}, \dots, x_{jp})$ represent p measurements for the j th individual. The distance between individuals i and j can be calculated using

Euclidean Distance: $\sqrt{\sum_{k=1}^p (x_{ki} - x_{kj})^2}$

Manhattan Distance: $\sum_{k=1}^p |x_{ki} - x_{kj}|$

Correlation Based Distance: $1 - P(x_i, x_j)$
 $P(x_i, x_j)$ can be Pearson's, Spearman's, or Kendall's correlations

- For categorical variables, suppose we have binary data as follows

		individual j	
		Present(+)	Absent(-)
individual i	Present(+)	a	b
	Absent(-)	c	d

Jaccard :

Similarity: $\frac{a}{ab+bc}$

distance: $1 - \frac{a}{ab+bc}$

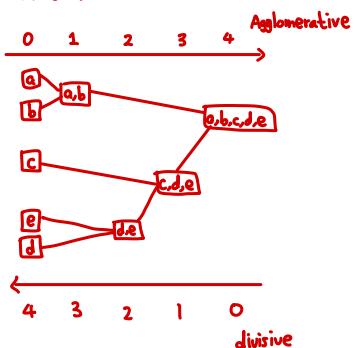
Simple Matching

Similarity: $\frac{a+d}{ab+cd}$

distance: $\frac{b+c}{ab+cd}$

Hierarchical Clustering

- 2 strategies: agglomerative or divisive
- Distances between individuals or clusters of individuals are calculated at each iteration



Distance matrices are calculated iteratively until all observations are grouped together in one cluster (agglomerative clustering) or until all individuals become separate clusters on their own (divisive clustering)

To calculate distances between clusters, the mechanism by which we calculate the distance between clusters at each iteration of the clustering algorithm is referred to as the linkage method

Linkage methods:

- Single Linkage: distance determined by considering the distance between the closest pair of individuals in the clusters
- Complete Linkage: distance between the farthest pair of individuals in the clusters

Average Linkage: distances between all pairs of individuals are first calculated, then the distances between clusters are then determined by taking all the distances. Clusters with the smallest average distance are joined together in the subsequent iteration

Gower Similarity

- Useful when there's a mix of categorical and continuous variables
- Gower's similarity first calculates similarity with respect to each of the variables, and takes the average similarity across the variables
- Let S_{ijk} be the similarity between individuals i and j , with respect to the k th variable. Gower's similarity between individuals i and j is then calculated as

$$S_{ij} = \frac{\sum w_{ijk} S_{ijk}}{\sum w_{ijk}}$$

$w_{ijk}=0$ when a given variable is missing for either of the individuals, $w_{ijk}=1$ otherwise

- To calculate S_{ijk} , then we have that

- If the k th variable is categorical, the simple matching similarity is used.
- If k th variable is continuous, then $S_{ijk} = 1 - \frac{|x_{ik} - x_{jk}|}{R_k}$, where x_{ik} and x_{jk} are the values of the variable corresponding to the i th and j th individuals, respectively. R_k represents the range of the k th variable
- Gower's distance: $d_{ij} = 1 - S_{ij}$

K-means Clustering

- Data first partitioned into k groups C_1, C_2, \dots, C_k
- Individuals then assigned to the k clusters, where class assignment is determined by using some optimization criteria
- Most commonly used criteria is the group (within cluster) sum of squares (WSS), where data is partitioned in such a way that the total WSS is minimized
- General idea: find all possible partitions, and choose the one that leads to minimum total WSS
- Not computationally feasible
- The algorithm:

- Select initial partition of data, using initially guessed number of clusters k
- Move each of the individuals around and calculate the change in the total WSS for all the moves made
- Accept the move that lead to the maximum improvement in the WSS - these are our newly established clusters
- Repeat 2 and 3 until no improvement is possible

Canonical Correlation Analysis

- Involves projections in direction of maximum correlation between two sets of variables
- Involves 2 data matrices X and Y
- Main objective is to discover relationships between the variables in the 2 datasets
- CCA seeks to find linear combinations of variables that sequentially maximize correlation between the two sets of variables

Covariance Matrices

- Consider 2 random vectors of length p and q

$$X = \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_p \end{pmatrix} \quad Y = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_q \end{pmatrix}$$

$\begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_p \end{pmatrix}$ combined vector

- The variance covariance matrix of the combined vector can be written as

$$\Sigma = \begin{pmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{pmatrix}$$

Σ_{11} covariance matrix of X

Σ_{22} covariance matrix of Y

$\Sigma_{12} = \Sigma_{21}$ covariances between X variables and Y variables

Σ_{11} and Σ_{22} within covariances

Σ_{12} and Σ_{21} cross covariances

CCA and Eigenvalues

- We can convert finding the linear combinations into a problem involving eigenvalues and eigenvectors

Rewrite the correlation as

$$\text{Cor}(z_i, w_i) = \frac{\underline{z}_i' \Sigma^{-\frac{1}{2}} \Sigma_{12} \Sigma_{22}^{-\frac{1}{2}} \underline{w}_i}{\sqrt{\underline{z}_i' \Sigma^{-\frac{1}{2}} \Sigma_{12} \Sigma_{22}^{-\frac{1}{2}} \underline{w}_i} \sqrt{\underline{w}_i' \Sigma^{-\frac{1}{2}} \Sigma_{22} \Sigma_{12}^{-\frac{1}{2}} \underline{z}_i}}$$

Let $\underline{a}' = \underline{y}_i' \Sigma^{-\frac{1}{2}}$ and $\underline{b}' = \Sigma_{22}^{-\frac{1}{2}} \underline{w}_i$. Then

$$\text{Cor}(z_i, w_i) = \frac{\underline{a}' \Sigma^{-\frac{1}{2}} \Sigma_{12} \Sigma_{22}^{-\frac{1}{2}} \underline{b}'}{\sqrt{\underline{a}' \Sigma^{-\frac{1}{2}}} \sqrt{\underline{b}' \Sigma^{-\frac{1}{2}}}}$$

under $\underline{a}' \underline{a} = 1$ and $\underline{b}' \underline{b} = 1$

Let $\Omega = \Sigma_{11}^{-\frac{1}{2}} \Sigma_{12} \Sigma_{22}^{-\frac{1}{2}}$. Then

$$\text{Cor}(z_i, w_i) = \frac{\underline{a}' \Omega \underline{b}'}{\sqrt{\underline{a}' \Omega} \sqrt{\underline{b}' \Omega}}$$

We can obtain

$$\text{Cor}(z_i, w_i)^2 = (\underline{a}' \Omega \underline{b}')^2 \leq \underline{a}' \Omega \Omega' \underline{b}'$$

$$\Omega \Omega' = \Sigma_{11}^{-\frac{1}{2}} \Sigma_{12} \Sigma_{22}^{-\frac{1}{2}} \Sigma_{22}^{-\frac{1}{2}} \Sigma_{12} \Sigma_{11}^{-\frac{1}{2}}$$

Let $\Psi = \Omega \Omega'$ and $\Psi = T \Lambda T'$ be its eigen decomposition. Then

$$\begin{aligned} \text{Cor}(z_i, w_i)^2 &\leq \underline{a}' \Omega \Omega' \underline{b}' = \underline{a}' T \Lambda T' \underline{b}' \\ &= \sum_{i=1}^n \lambda_i a_i^2 \leq \lambda_1 \sum_{i=1}^n a_i^2 = \lambda_1 \end{aligned}$$

λ_1 eigenvalues of $\Omega \Omega'$

λ_1 maximum eigenvalue

Then maximum correlation is $\sqrt{\lambda_1}$ and the standardized eigenvector corresponding to λ_1 satisfies $\underline{a}' \underline{a} = 1$

$$\underline{a}' = \underline{y}_i' \Sigma^{-\frac{1}{2}} \Rightarrow \underline{y}_i' = \underline{a}' \Sigma^{-\frac{1}{2}}$$

Similarly, we can get

$$\text{Cor}(w_i, z_i) = \left(\frac{\underline{a}' \Omega \underline{b}'}{\sqrt{\underline{a}' \Omega} \sqrt{\underline{b}' \Omega}} \right)^2 \leq \underline{b}' \Omega' \Omega \underline{a}$$

Leads to same eigenvalues but different eigenvectors

Choose \underline{b} standardized eigenvector of $\Omega' \Omega$ corresponding to λ_1 ,

$$\underline{b} = \Sigma_{22}^{-\frac{1}{2}} \underline{w}_i \Rightarrow \underline{w}_i = \Sigma_{22}^{-\frac{1}{2}} \underline{b}$$

To get k th pair of canonical variates, we use the standardized eigenvectors corresponding to the k th largest eigenvalues of $\Omega \Omega'$ and $\Omega' \Omega$

Correlations sequentially maximized

$$\sqrt{\lambda_1} \geq \sqrt{\lambda_2} \geq \dots \geq \sqrt{\lambda_n}$$

We can obtain that

$\Omega \Omega'$ and $\Sigma_{11}^{-\frac{1}{2}} \Sigma_{12} \Sigma_{22}^{-\frac{1}{2}} \Sigma_{22}^{-\frac{1}{2}}$ are similar

$\Omega' \Omega$ and $\Sigma_{22}^{-\frac{1}{2}} \Sigma_{12} \Sigma_{11}^{-\frac{1}{2}} \Sigma_{12} \Sigma_{22}^{-\frac{1}{2}}$ are similar

Then we obtain

$$\underline{a}' = \Sigma^{-\frac{1}{2}} \underline{a} \quad \underline{a}' \text{ eigenvector of } \Sigma_{11}^{-\frac{1}{2}} \Sigma_{12} \Sigma_{22}^{-\frac{1}{2}}$$

$$\underline{b}' = \Sigma_{22}^{-\frac{1}{2}} \underline{b} \quad \underline{b}' \text{ eigenvector of } \Sigma_{22}^{-\frac{1}{2}} \Sigma_{12} \Sigma_{11}^{-\frac{1}{2}} \Sigma_{12} \Sigma_{22}^{-\frac{1}{2}}$$

We can then get

$$\underline{a}' = \underline{a}'' \quad \underline{a}'' \text{ eigenvector of } \Sigma_{11}^{-\frac{1}{2}}$$

$$\underline{b}' = \underline{b}'' \quad \underline{b}'' \text{ eigenvector of } \Sigma_{22}^{-\frac{1}{2}}$$

CCA and SVD

We can use the singular value decomposition of $\Omega = \Sigma_{11}^{-\frac{1}{2}} \Sigma_{12} \Sigma_{22}^{-\frac{1}{2}} = PSQ'$

P and Q orthogonal matrices

S diagonal matrix with singular values as its diagonal elements

Then we have the solutions are

$$V = \Sigma_{11}^{-\frac{1}{2}} P \quad U = \Sigma_{22}^{-\frac{1}{2}} Q$$

where $V = (v_1 \ v_2 \ \dots \ v_r)$ and $U = (u_1 \ u_2 \ \dots \ u_r)$

Uncorrelatedness of Canonical Variates

- We can easily show the canonical variates ξ and η are uncorrelated and have unit variance

$$\begin{aligned}\text{Cov}(\xi) &= \text{Cov}(V'\boldsymbol{\lambda}) \\ &= V' \text{Cov}(\boldsymbol{\lambda}) V \\ &= P' \Sigma_{\boldsymbol{\lambda}}^{-\frac{1}{2}} \Sigma_{\boldsymbol{\lambda}} \Sigma_{\boldsymbol{\lambda}}^{-\frac{1}{2}} P \\ &= P' P = I\end{aligned}$$

$$\begin{aligned}\text{Cov}(\eta) &= \text{Cov}(U'\boldsymbol{\gamma}) \\ &= U' \text{Cov}(\boldsymbol{\gamma}) U \\ &= Q' \Sigma_{\boldsymbol{\gamma}}^{-\frac{1}{2}} \Sigma_{\boldsymbol{\gamma}} \Sigma_{\boldsymbol{\gamma}}^{-\frac{1}{2}} Q \\ &= Q' Q = I\end{aligned}$$

Some more Remarks on CCA

- We often use correlation matrices to conduct CCA because variables may be on different scales
- We cannot obtain the population covariance matrices in practice, so we can use the sample covariance matrices

$$\hat{\Sigma} = \begin{pmatrix} \hat{\Sigma}_{xx} & \hat{\Sigma}_{xy} \\ \hat{\Sigma}_{yx} & \hat{\Sigma}_{yy} \end{pmatrix} = \begin{pmatrix} S_{xx} & S_{xy} \\ S_{yx} & S_{yy} \end{pmatrix}$$

Presenting CCA Results

- The main CCA results are the canonical variates and the corresponding correlation between the pairs of canonical variates
- We also explore and interpret the contributions of the original set of variables to the corresponding canonical variates
- One way to present is using finger plots



- Canonical coefficients can also be plotted with bar graphs or a helio-plot