

STAT 475/575 Midterm Review

Outline

Data

- Vector-format: $\mathbf{X} = \begin{bmatrix} x_1 \\ \vdots \\ x_p \end{bmatrix}$
- Matrix-format (n observations): $X = \begin{bmatrix} x_{11} & \cdots & x_{1p} \\ \vdots & & \vdots \\ x_{n1} & \cdots & x_{np} \end{bmatrix}$

Multivariate normal distribution (MVN)

- $\mathbf{X} \sim N_p(\boldsymbol{\mu}, \Sigma)$
- PDF: $f(\mathbf{x}) = \frac{1}{(2\pi)^{p/2} |\Sigma|^{1/2}} \exp \left[-\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})' \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu}) \right]$
- Contour is ellipse/ellipsoid: $(\mathbf{x} - \boldsymbol{\mu})' \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu}) = c$
- Check normality:
 - Plot: histogram, QQ plot, ...
 - Test: Shapiro-Wilk test, ...

Compare center (hypothesis test)

- $H_0 : \dots = \mathbf{c}$ vs $H_a : \dots \neq \mathbf{c}$
- Assumption/condition
- Test statistics & distribution
- Case
 - One-sample (Hotelling's T²): $\boldsymbol{\mu} = \mathbf{c}$
 - Two-sample (Hotelling's T²): $\boldsymbol{\mu}^{(1)} - \boldsymbol{\mu}^{(2)} = \mathbf{c}$
 - More samples (MANOVA): $\boldsymbol{\mu}^{(1)} = \dots = \boldsymbol{\mu}^{(g)}$
- Confidence region
 - Region/ellipsoid (p -dim)
 - Interval (1-dim)
 - * One-at-a-time
 - * Bonferroni (simultaneous)
 - * T^2 (simultaneous)

Principal Component Analysis (PCA)

$$PC_j = \sum_{i=1}^p a_{ji} X_i$$

Factor Analysis (FA)

$$X_i - \mu_i = \sum_{j=1}^m l_{ij} F_j + \epsilon_i.$$

- Principal component method
- Principal factor method
- Maximum likelihood estimation

One-sample Hotelling's T2 Test

Univariate case: one-sample t -test

$H_0 : \mu = \mu_0$ vs $H_a : \mu \neq \mu_0$.

$$t = \frac{\bar{X} - \mu_0}{s/\sqrt{n}} \sim t_{n-1}$$

Multivariate case: one-sample Hotelling's T^2 test

$H_0 : \boldsymbol{\mu} = \boldsymbol{\mu}_0$ vs $H_a : \boldsymbol{\mu} \neq \boldsymbol{\mu}_0$ (at least one variable not equal).

$$T^2 = (\bar{\mathbf{X}} - \boldsymbol{\mu}_0)' \left(\frac{\mathbf{S}}{n} \right)^{-1} (\bar{\mathbf{X}} - \boldsymbol{\mu}_0) \sim \frac{(n-1)p}{n-p} F_{p, n-p}$$

where sample covariance $\mathbf{S} = \frac{1}{n-1} \sum_i (\mathbf{X}_i - \bar{\mathbf{X}})(\mathbf{X}_i - \bar{\mathbf{X}})'$.

In R: `DescTools::HotellingsT2Test`

Conditions

1. Independency (most important)
2. Sample from a population with mean $\boldsymbol{\mu}$, covariance Σ
3. Multivariate-normality (in R: `mvShapiroTest::mvShapiro.Test`)

Repeated measure

$H_0 : C\boldsymbol{\mu} = \mathbf{c}$ vs $H_a : C\boldsymbol{\mu} \neq \mathbf{c}$.

$$T^2 = (C\bar{\mathbf{X}} - \mathbf{c})' \left(\frac{C\mathbf{S}C'}{n} \right)^{-1} (C\bar{\mathbf{X}} - \mathbf{c}) \sim \frac{(n-1)q}{n-q} F_{q, n-q}$$

where $q = \text{rank}(C)$.

Two-sample Hotelling's T2 Test

Univariate case: two-sample t -test

$H_0 : \mu_1 - \mu_2 = c$ vs $H_a : \mu_1 - \mu_2 \neq c$.

$$t = \frac{\bar{X}_1 - \bar{X}_2 - c}{\sqrt{s_1^2/n_1 + s_2^2/n_2}} \sim t_{n_1+n_2-2}$$

where s_1^2, s_2^2 can be replaced with pooled variance $s_p^2 = \frac{(n_1-1)s_1^2 + (n_2-1)s_2^2}{(n_1-1) + (n_2-1)}$.

Multivariate case: two-sample Hotelling's T^2 test

$H_0 : \boldsymbol{\mu}^{(1)} - \boldsymbol{\mu}^{(2)} = \mathbf{c}$ vs $H_a : \boldsymbol{\mu}^{(1)} - \boldsymbol{\mu}^{(2)} \neq \mathbf{c}$ (at least one variable not equal).

$$T^2 = (\bar{\mathbf{X}}^{(1)} - \bar{\mathbf{X}}^{(2)} - \mathbf{c})' \left(\frac{\mathbf{S}_1}{n_1} + \frac{\mathbf{S}_2}{n_2} \right)^{-1} (\bar{\mathbf{X}}^{(1)} - \bar{\mathbf{X}}^{(2)} - \mathbf{c}) \sim \frac{(n_1 + n_2 - 2)p}{n_1 + n_2 - p - 1} F_{p, n_1 + n_2 - p - 1}$$

where $\mathbf{S}_1, \mathbf{S}_2$ have to be replaced with pooled variance $\mathbf{S}_p = \frac{(n_1-1)\mathbf{S}_1 + (n_2-1)\mathbf{S}_2}{(n_1-1) + (n_2-1)}$.

In R: `DescTools::HotellingsT2Test`

Conditions

1. Independency (within group & between group) (most important)
2. Sample 1 from a population with mean $\boldsymbol{\mu}_1$, covariance Σ_1
Sample 2 from a population with mean $\boldsymbol{\mu}_2$, covariance Σ_2
3. Homogeneity of variance ($\Sigma_1 = \Sigma_2$) (in R: `biotools::boxM`)
4. Multivariate-normality (for each group) (in R: `mvShapiroTest::mvShapiro.Test`)

MANOVA ($g > 2$)

Univariate case: ANOVA

$H_0 : \mu_1 = \dots = \mu_g$ vs $H_a : \mu_k \neq \mu_l$ for some k, l .

$$F = \frac{MS_{trt}}{MS_{err}} = \frac{SS_{trt}/(g-1)}{SS_{err}/(n-g)} \sim F_{g-1, n-g}$$

Multivariate case: MANOVA

$H_0 : \boldsymbol{\mu}^{(1)} = \dots = \boldsymbol{\mu}^{(g)}$ vs $H_a : \boldsymbol{\mu}^{(k)} \neq \boldsymbol{\mu}^{(l)}$ for some k, l (at least two groups, at least one variable not equal).

$$(\text{Wilk's}) \Lambda = \frac{|W|}{|B+W|} = \frac{\left| \sum_{l=1}^g \sum_{j=1}^{n_l} (\mathbf{X}_j^{(l)} - \bar{\mathbf{X}}^{(l)})(\mathbf{X}_j^{(l)} - \bar{\mathbf{X}}^{(l)})' \right|}{\left| \sum_{l=1}^g \sum_{j=1}^{n_l} (\mathbf{X}_j^{(l)} - \bar{\mathbf{X}})(\mathbf{X}_j^{(l)} - \bar{\mathbf{X}})' \right|} \quad i.e. \frac{|\text{Within-group SS}|}{|\text{Total SS}|}$$

In R: `car::Manova`

Conditions

1. Independency (within group & between group) (most important)
2. Each sample from a population with mean $\boldsymbol{\mu}_l$, covariance Σ_l
3. Homogeneity of variance ($\Sigma_1 = \dots = \Sigma_g$) (in R: `biotools::boxM`)
4. Multivariate-normality (for each group) (in R: `mvShapiroTest::mvShapiro.Test`)

Confidence Region/Interval

(p -dim) Region for $\mu, \mu^{(1)} - \mu^{(2)}$

One-sample:

$$T^2 = (\bar{X} - \mu_0)' \left(\frac{S}{n} \right)^{-1} (\bar{X} - \mu_0) \sim \frac{(n-1)p}{n-p} F_{p, n-p}$$

$$\Rightarrow \mu : (\bar{X} - \mu)' \left(\frac{S}{n} \right)^{-1} (\bar{X} - \mu) \leq \frac{(n-1)p}{n-p} F_{(p, n-p), 1-\alpha}$$

Two-sample:

$$T^2 = (\bar{X}^{(1)} - \bar{X}^{(2)} - c)' \left(\frac{S_1}{n_1} + \frac{S_2}{n_2} \right)^{-1} (\bar{X}^{(1)} - \bar{X}^{(2)} - c) \sim \frac{(n_1 + n_2 - 2)p}{n_1 + n_2 - p - 1} F_{p, n_1 + n_2 - p - 1}$$

$$\Rightarrow \mu^{(1)} - \mu^{(2)} : (\bar{X}^{(1)} - \bar{X}^{(2)} - (\mu^{(1)} - \mu^{(2)}))' \left(\frac{S_1}{n_1} + \frac{S_2}{n_2} \right)^{-1} (\bar{X}^{(1)} - \bar{X}^{(2)} - (\mu^{(1)} - \mu^{(2)}))$$

$$\leq \frac{(n_1 + n_2 - 2)p}{n_1 + n_2 - p - 1} F_{(p, n_1 + n_2 - p - 1), 1-\alpha}$$

where S_1, S_2 have to be replaced with pooled variance S_p .

(1-dim) Interval for $\mu_j, \mu_j^{(1)} - \mu_j^{(2)}, \mu_j^{(k)} - \mu_j^{(l)}$

Formula:

$$\text{center} \pm \text{coef} \times \text{std.err}$$

If $g \geq 2$:

- individual variances: $S_{jj}^{(l)}$
- pooled variance: $S_{p,jj} = \frac{1}{n-g} \sum_{l=1}^g (n_l - 1) S_{jj}^{(l)}$

(1). One-at-a-time CI

$$\mu_j : \left[\bar{X}_j \pm t_{n-1, 1-\alpha/2} \cdot \sqrt{\frac{S_{jj}}{n}} \right]$$

$$\mu_j^{(1)} - \mu_j^{(2)} : \left[\bar{X}_j^{(1)} - \bar{X}_j^{(2)} \pm t_{n_1+n_2-2, 1-\alpha/2} \cdot \sqrt{\frac{S_{jj}^{(1)}}{n_1} + \frac{S_{jj}^{(2)}}{n_2}} \right]$$

$$\mu_j^{(k)} - \mu_j^{(l)} : \left[\bar{X}_j^{(k)} - \bar{X}_j^{(l)} \pm t_{n-g, 1-\alpha/2} \cdot \sqrt{\frac{S_{jj}^{(k)}}{n_k} + \frac{S_{jj}^{(l)}}{n_l}} \right]$$

(2). Bonferroni simultaneous CI

Trick: \cdot/m in one-at-a-time CI, where m is the no. of comparisons (CIs) you want. E.g.

- (One sample) $\mu_1, \dots, \mu_p \Rightarrow m = p$
- (Two samples) $\mu_1^{(1)} - \mu_1^{(2)}, \dots, \mu_p^{(1)} - \mu_p^{(2)} \Rightarrow m = p$
- (More samples) $\begin{cases} \mu_1^{(1)} - \mu_1^{(2)}, \dots, \mu_p^{(1)} - \mu_p^{(2)} \\ \mu_1^{(1)} - \mu_1^{(3)}, \dots, \mu_p^{(1)} - \mu_p^{(3)} \\ \mu_1^{(2)} - \mu_1^{(3)}, \dots, \mu_p^{(2)} - \mu_p^{(3)} \end{cases} \Rightarrow m = p \cdot \binom{g}{2}$

$$\mu_j : \left[\bar{X}_j \pm t_{n-1, 1-\alpha/2m} \cdot \sqrt{\frac{S_{jj}}{n}} \right]$$

$$\mu_j^{(1)} - \mu_j^{(2)} : \left[\bar{X}_j^{(1)} - \bar{X}_j^{(2)} \pm t_{n_1+n_2-2, 1-\alpha/2m} \cdot \sqrt{\frac{S_{jj}^{(1)}}{n_1} + \frac{S_{jj}^{(2)}}{n_2}} \right]$$

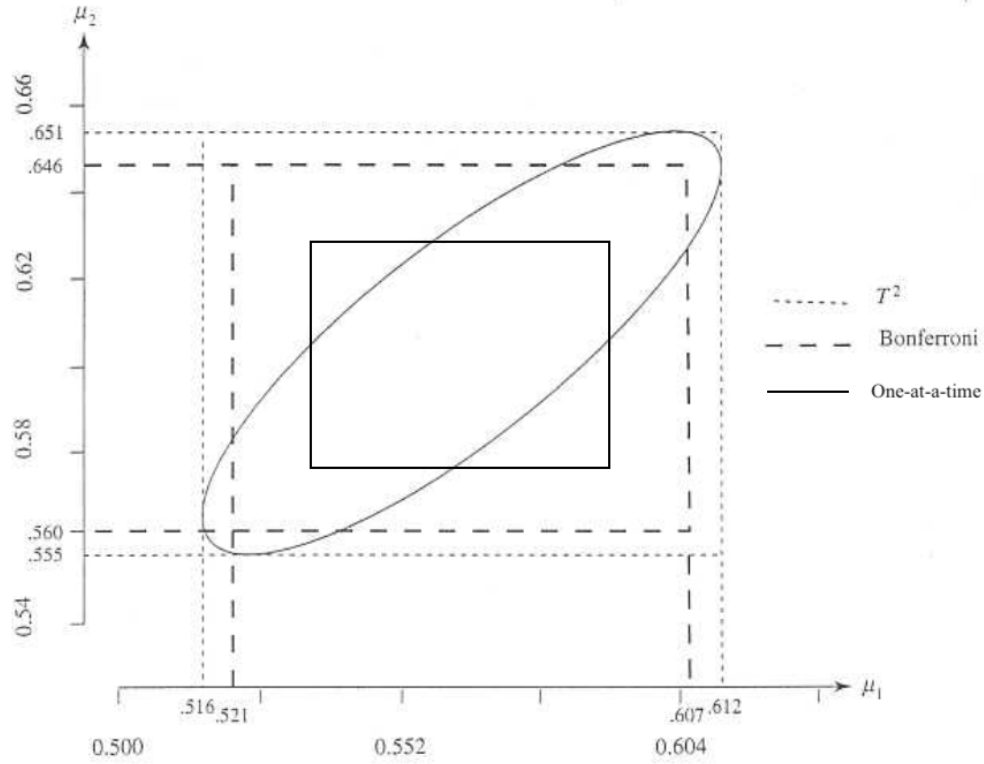
$$\mu_j^{(k)} - \mu_j^{(l)} : \left[\bar{X}_j^{(k)} - \bar{X}_j^{(l)} \pm t_{n-g, 1-\alpha/2m} \cdot \sqrt{\frac{S_{jj}^{(k)}}{n_k} + \frac{S_{jj}^{(l)}}{n_l}} \right]$$

(3). T^2 simultaneous CI

Trick: $\sqrt{\cdot}$ in CR.

$$\mu_j : \left[\bar{X}_j \pm \sqrt{\frac{(n-1)p}{n-p} F_{(p, n-p), 1-\alpha}} \cdot \sqrt{\frac{S_{jj}}{n}} \right]$$

$$\mu_j^{(1)} - \mu_j^{(2)} : \left[\bar{X}_j^{(1)} - \bar{X}_j^{(2)} \pm \sqrt{\frac{(n_1+n_2-2)p}{n_1+n_2-p-1} F_{(p, n_1+n_2-p-1), 1-\alpha}} \cdot \sqrt{\frac{S_{jj}^{(1)}}{n_1} + \frac{S_{jj}^{(2)}}{n_2}} \right]$$



Principal Component Analysis (PCA)

Model

Idea: uncorrelated linear combinations of original covariates \mathbf{X} .

$$\mathbf{Y}_{p \times 1} = \mathbf{A}'_{p \times p} \mathbf{X}_{p \times 1} \quad \begin{cases} Y_1 = \mathbf{a}'_1 \mathbf{X} = a_{11}X_1 + a_{12}X_2 + \cdots + a_{1p}X_p \\ Y_2 = \mathbf{a}'_2 \mathbf{X} = a_{21}X_1 + a_{22}X_2 + \cdots + a_{2p}X_p \\ \vdots \\ Y_p = \mathbf{a}'_p \mathbf{X} = a_{p1}X_1 + a_{p2}X_2 + \cdots + a_{pp}X_p \end{cases}$$

$Var(\mathbf{X}) = \Sigma$ with eigenvalues $\lambda_1 \geq \cdots \geq \lambda_p \geq 0$. $\mathbf{a}'_k \mathbf{a}_k = 1$, $Cov(Y_k, Y_l) = \mathbf{a}'_k \Sigma \mathbf{a}_l = 0$.

Properties

- Spectral(eigen) decomposition: $\Sigma = E \Lambda E' = \lambda_1 \mathbf{e}_1 \mathbf{e}'_1 + \cdots + \lambda_p \mathbf{e}_p \mathbf{e}'_p$, with eigenvalue-eigenvector pairs $(\lambda_i, \mathbf{e}_i)$.
- Loading matrix: $A = E = [\mathbf{e}_1, \cdots, \mathbf{e}_p]$, with coefficients $\mathbf{a}_j = \mathbf{e}_j$.
- $Var(\mathbf{Y}) = \Lambda = \text{diag}(\lambda_i)$ and $Cov(\mathbf{Y}, \mathbf{X}) = \Lambda E'$.
 $Var(Y_j) = \lambda_j$ or $\frac{\lambda_j}{\lambda_1 + \cdots + \lambda_p}$ (proportion).
- PC score: $\mathbf{Y}_{n \times p} = \mathbf{X}_{n \times p} A$ or $\mathbf{Y}_{p \times 1} = \mathbf{A}' \mathbf{X}$.
- Choose number of PCs: (1) scree plot + elbow method (2) % variance (3) context (4) meaningful interpretation (5) desire for simplicity (6)...
- Interpretation: size (absolute value) and sign (positive or negative) of coefficients \mathbf{a}_j .

In R: `prcomp` or `princomp`

`prcomp(x, center = T, scale. = F)`

- `scale. = F`: raw/centered data (covariance)
- `scale. = T`: standardized data (correlation)

`princomp(x, cor = F)` or `princomp(covmat, cor = F)`

- `cor = F`: centered data (covariance)
- `cor = T`: standardized data (correlation)

Factor Analysis (FA)

Model

Idea: latent variables behind original covariates \mathbf{X} .

$$\mathbf{X}_{p \times 1} - \boldsymbol{\mu}_{p \times 1} = \mathbf{L}_{p \times m} \mathbf{F}_{m \times 1} + \boldsymbol{\epsilon}_{p \times 1} \quad \begin{cases} X_1 - \mu_1 = l_{11}F_1 + l_{12}F_2 + \dots + l_{1m}F_m + \epsilon_1 \\ X_2 - \mu_2 = l_{21}F_1 + l_{22}F_2 + \dots + l_{2m}F_m + \epsilon_2 \\ \vdots \\ X_p - \mu_p = l_{p1}F_1 + l_{p2}F_2 + \dots + l_{pm}F_m + \epsilon_p \end{cases}$$

$$E(\mathbf{F}) = \mathbf{0}, \text{Var}(\mathbf{F}) = \mathbf{I}_m, E(\boldsymbol{\epsilon}) = \mathbf{0}, \text{Var}(\boldsymbol{\epsilon}) = \Psi = \text{diag}(\psi_i), \text{ and } \mathbf{F} \perp \boldsymbol{\epsilon}.$$

Properties

- $\text{Var}(\mathbf{X}) = \Sigma = \mathbf{L}\mathbf{L}' + \Psi$ and $\text{Cov}(\mathbf{X}, \mathbf{F}) = \mathbf{L}$.
 $\text{Var}(X_i) = \sigma_{ii} = l_{i1}^2 + \dots + l_{im}^2 + \psi_i$ and $\text{Cov}(X_i, X_k) = \sigma_{ik} = l_{i1}l_{k1} + \dots + l_{im}l_{km}$.
- $\sigma_{ii} = h_i^2 + \psi_i$. Communality: $h_i^2 = \sum_{j=1}^m l_{ij}^2$. Uniqueness: $\psi_i = \sigma_{ii} - h_i^2$.
- Infinite rotations: $\mathbf{X} - \boldsymbol{\mu} = \mathbf{L}\mathbf{F} + \boldsymbol{\epsilon} = (\mathbf{L}\mathbf{T})(\mathbf{T}'\mathbf{F}) + \boldsymbol{\epsilon}$ for any orthogonal matrix \mathbf{T} .
 E.g. varimax, quartimax, promax, ...
- Contribution of j -th factor: $\sum_{i=1}^p l_{ij}^2$ or $\frac{\sum_{i=1}^p l_{ij}^2}{\lambda_1 + \dots + \lambda_p}$ (proportion).
- Factor score: $\mathbf{F}_{n \times m} = (\mathbf{X}_{n \times p} - \mathbf{1}_n \bar{\mathbf{X}}')(\mathbf{L}\mathbf{L}' + \Psi)^{-1}\mathbf{L}$ or $\mathbf{F}_{m \times 1} = \mathbf{L}'(\mathbf{L}\mathbf{L}' + \Psi)^{-1}(\mathbf{X} - \bar{\mathbf{X}})$.
- Interpretation: size (absolute value) + sign (positive or negative) of coefficients l_j .

Three common methods of estimating factor loadings:

1. Principal component method:

- Loading matrix: $\mathbf{L} = [\mathbf{l}_1, \dots, \mathbf{l}_m] = [\sqrt{\lambda_1}\mathbf{e}_1, \dots, \sqrt{\lambda_m}\mathbf{e}_m]$, with coefficients $l_j = \sqrt{\lambda_j}e_j$.
- Choose number of Factors: (1) no. of PCs in PCA (2) % variance (3) context (4) meaningful interpretation (5) desire for simplicity (6)...

2. Principal factor method:

- Iterative method: $\Psi \rightarrow \mathbf{L} \rightarrow \Psi \rightarrow \mathbf{L} \rightarrow \dots$.

3. Maximum likelihood estimation:

- Normality assumption: $\mathbf{X}, \mathbf{F}, \boldsymbol{\epsilon}$ are normal.
- Maximize the likelihood function with respect to $(\boldsymbol{\mu}, \mathbf{L}, \Psi)$.
- Likelihood ratio test: H_0 : m -factors are sufficient.
- Choose number of Factors: (1) likelihood ratio test (2) % variance (3) context (4) meaningful interpretation (5) desire for simplicity (6)...

In R: `prcomp` or `factanal`

1. PC method:

- (1) PCA: `prcomp` or `princomp`
- (2) Factor loading: `rotation %*% diag(sdev)`
- (3) Rotation (optional): `varimax` on first m columns of factor loading

3. MLE:

`factanal(x, factors, scores = "regression", rotation = "varimax")` or
`factanal(factors, covmat, n.obs, scores = "regression", rotation = "varimax")`

- `rotation`: "varimax" (default) or "none"
- `scores`: "none" (default) or "regression"