

Lecture 10 (Mar 19) CANONICAL CORR. ANALYSIS

$$\underline{X} \in \mathbb{R}^{n \times p}$$

$$\underline{Y} \in \mathbb{R}^{n \times q}$$

↓

$$X \in \mathbb{R}^p$$

↓

$$Y \in \mathbb{R}^q$$

$$(\underline{X}, \underline{Y}) \in \mathbb{R}^{n \times (p+q)}$$

$$\downarrow$$
$$(X, Y)$$

POPULATION CASE $X \in \mathbb{R}^p, Y \in \mathbb{R}^q$

$$\Sigma_{XX} = \text{Var}(X) \quad \Sigma_{YY} = \text{Var}(Y)$$

$$\Sigma_{XY} = \text{Cov}(X, Y)$$

GOAL: Find $a \in \mathbb{R}^p, b \in \mathbb{R}^q$

$$\boxed{\text{corr}(a^T X, b^T Y) \rightarrow \max}$$

$$g(a, b) = \text{corr}(a^T X, b^T Y)$$

$$= \frac{\text{cov}(a^T X, b^T Y)}{\sqrt{\text{Var}(a^T X) \text{Var}(b^T Y)}}$$

$$\sqrt{\text{Var}(a^T X) \text{Var}(b^T Y)}$$

$$= \frac{\sqrt{\text{var}(a \cdot x) \cdot \text{var}(b \cdot y)}}{a^T \Sigma_{xy} b}$$

$$= \frac{a^T \Sigma_{xx} a \cdot b^T \Sigma_{yy} b}{\sqrt{a^T \Sigma_{xx} a \cdot b^T \Sigma_{yy} b}}$$

$$\alpha := \Sigma_{xx}^{1/2} a \quad \alpha^T \alpha = a^T \Sigma_{xx} a$$

$$\beta := \Sigma_{yy}^{1/2} b \quad \beta^T \beta = b^T \Sigma_{yy} b$$

$$a = \Sigma_{xx}^{-1/2} \alpha \quad b = \Sigma_{yy}^{-1/2} \beta$$

$$= \frac{\alpha^T \Sigma_{xx}^{-1/2} \Sigma_{xy} \Sigma_{yy}^{-1/2} \beta}{\sqrt{\alpha^T \alpha} \cdot \sqrt{\beta^T \beta}}$$

$$= \left(\frac{\alpha}{\|\alpha\|} \right)^T M \left(\frac{\beta}{\|\beta\|} \right) \in \mathbb{R}^{p \times q}$$

$$M = \Sigma_{xx}^{-1/2} \Sigma_{xy} \Sigma_{yy}^{-1/2}$$

Solve

maximize $\alpha^T M \beta$ s.t

$$\|\alpha\| = \|\beta\| = 1$$

Lagrangian

$$2\alpha^T M \beta - \sigma(\alpha^T \alpha - 1) - \sigma'(\beta^T \beta - 1)$$

$$\nabla_{\alpha} = 2 \cdot M \beta - 2\sigma \alpha$$

$$\nabla_{\beta} = 2M^T \alpha - 2\sigma' \beta$$

FOC:

$$M \beta = \sigma \alpha$$

$$M^T \alpha = \sigma' \beta$$

$$\alpha^T M \beta = \sigma \cdot \alpha^T \alpha = \sigma$$

$$\beta^T M^T \alpha = \sigma' \beta^T \beta = \sigma'$$

$$\boxed{\sigma = \sigma'}$$

$$M\beta = \sigma\alpha$$

$$M^T\alpha = \sigma\beta$$

so...

$$M^T M \beta = \sigma \cdot M^T \alpha = \sigma^2 \beta$$

$$M M^T \alpha = \sigma M \beta = \sigma^2 \alpha$$

so

β is eigenv. of $M^T M$
with eigenvalue σ^2

α is eigenv. of $M M^T$
with eigenvalue σ^2

we maximize

$$\alpha^T M \beta = \sigma$$

so pick max eigenv.

Note: $M = U D V^T$

$$M^T M = V D^T D V^T$$

$$M M^T = U D D^T U^T$$

So α is the 1st col. of U
 β is the 1st col of V .

Q: Say we got (α_1, β_1)

$$a_1 = \sum x x^{-1/2} \alpha_1$$

$$b_1 = \Sigma_{YY}^{-1/2} \beta_1$$

Find α, β s.t

$$\underline{\alpha^\top \alpha_1 = 0}$$

||

$$\alpha^\top \Sigma_{XX} \alpha_1 = \text{Cov}(\alpha^\top X, \alpha_1^\top X)$$

$$\underline{\beta^\top \beta_1 = 0}$$

||

$$\text{Cov}(b^\top Y, b_1^\top Y)$$

so that

$$\text{Covr}(\alpha^\top X, b^\top Y) \rightarrow \max$$

$$\text{Covr}(\alpha^\top \Sigma_{XX}^{-1/2} X, \beta^\top \Sigma_{YY}^{-1} Y)$$

Solution:

$$\hat{\alpha} = 2^{\text{nd}} \text{ col. of } U$$

$$\hat{\beta} = 2^{\text{nd}} \text{ col. of } V$$

$$\eta_1 = a_1^T X$$

$$\phi_1 = b_1^T Y$$

$$\eta_2 = a_2^T X$$

$$\phi_2 = b_2^T Y$$

⋮

⋮

$$\eta_r = a_r^T X$$

$$\phi_r = b_r^T Y$$

$$\text{cov}(\eta_1, \phi_1) =$$

$$= a_1^T \Sigma_{xy} b_1$$

$$= a_1^T \underbrace{\left[\Sigma_{xx}^{-1/2} \Sigma_{xy} \Sigma_{yy}^{-1/2} \right]}_M \cdot \beta_1$$

$$= \sigma_1$$

$$\text{cov}(\eta_1, \phi_2) = \underbrace{\alpha_1^T M \beta_2^T}_{\text{2nd col. } V}$$

$$\left\{ \begin{array}{l} U D V^T \cdot \beta_2 = \sigma_2 \cdot \alpha_2 \end{array} \right.$$

$$= \sigma_2 \cdot \alpha_1^T \cdot \alpha_2 = 0$$

$$\text{cov}(\eta_1, \eta_1) \stackrel{a_1^T X}{=} a_1^T X$$

$$= a_1^T \Sigma_{XX} a_1$$

$$\approx a_1^T a_1 = 1$$

$$\text{Var} \left(\begin{array}{c} \eta_1 \\ \vdots \\ \eta_r \\ \phi_1 \\ \vdots \\ \phi_s \end{array} \right) = \left(\begin{array}{c|c} I_r & D_r \\ \hline D_r & I_r \end{array} \right)$$

$$D_r = \begin{pmatrix} \sigma_1 & & \\ & \ddots & \\ & & \sigma_r \end{pmatrix}$$

DATA:

$(x_i, y_i) \quad i = 1, \dots, n$

$n \geq p, q$

Compute the
sample covariance
matrix

$$S = \begin{pmatrix} S_{xx} & S_{xy} \\ S_{yx} & S_{yy} \end{pmatrix}$$

$$\hat{M} = S_{xx}^{-1/2} S_{xy} S_{yy}^{-1/2}$$

\wedge

$\wedge \rightarrow \wedge^T$

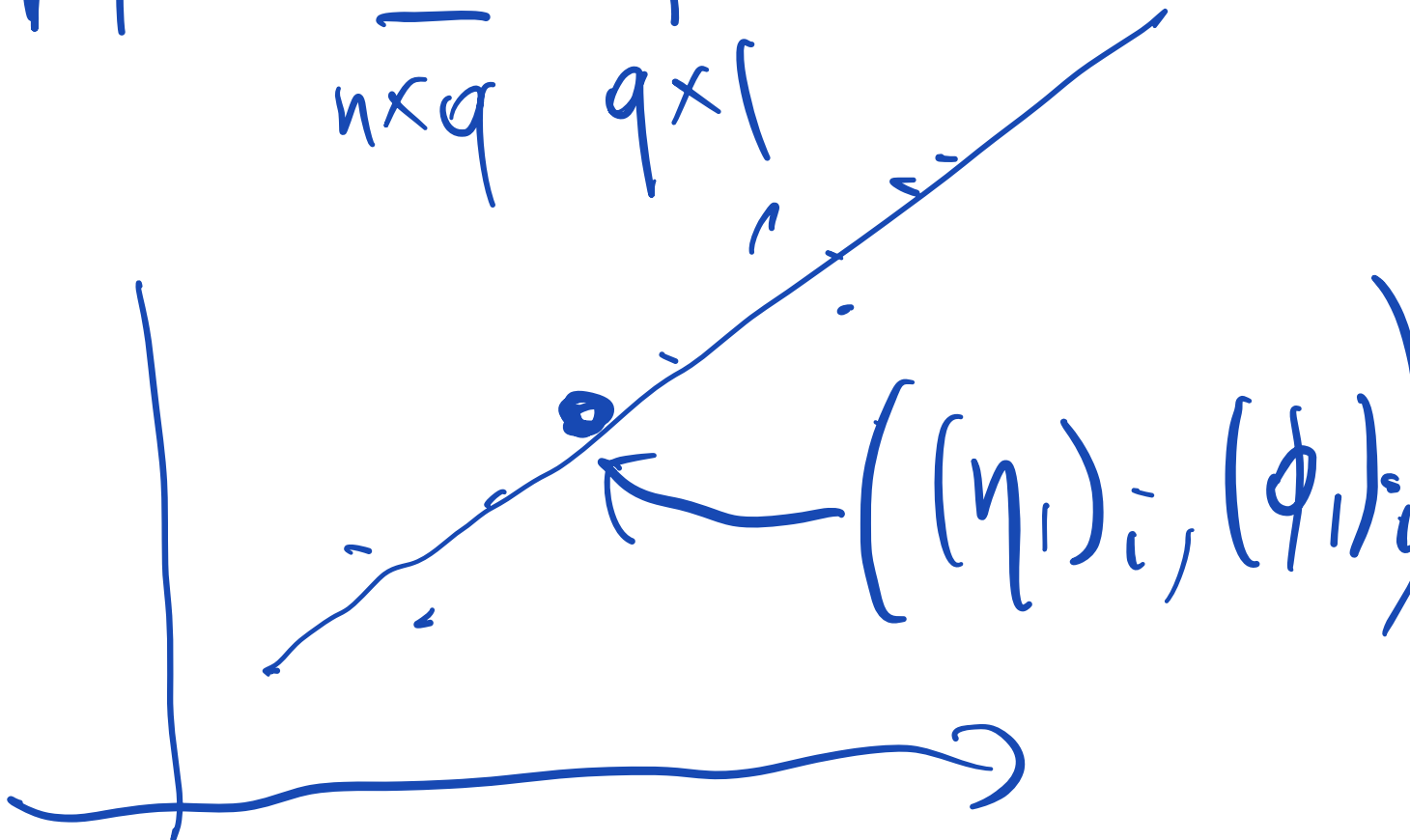
$$M = U D V$$

VISUALIZATION

$$\eta_i = \underbrace{X}_{n \times p} \cdot \underbrace{a_i}_{p \times 1}$$

biplot

$$\phi_i = \underbrace{Y}_{n \times q} \cdot \underbrace{b_i}_{q \times 1}$$



Heatmap

$$y_i = \underline{X} \cdot a_i$$

$$= \begin{pmatrix} a_i^T \underline{x}_1 \\ \vdots \\ a_i^T \underline{x}_n \end{pmatrix}$$

$$a_i^T \underline{x}_i \mapsto \mathbb{R}^D$$

STA 437/2005:
Methods for Multivariate Data
Week 10: Factor Analysis

Piotr Zwiernik

University of Toronto

Low-dimensional structures in multivariate statistics

Discussing covariance matrix estimation we mentioned some special structures that are routinely assumed in multivariate statistics.

We will discuss in detail two such constraints:

- ▶ $\Sigma = L + S$ where L is low-rank and S is sparse.
- ▶ the inverse of Σ^{-1} is sparse.

rank $\leq r$

In Factor Analysis: $\Sigma = WW^T + \Psi$ with $W \in \mathbb{R}^{m \times r}$, Ψ diagonal.

We explain here how such structure can occur by discussing motivating examples.

Factor Analysis: Motivating Examples

Example: Capital Asset Pricing Model (CAPM)

Models stock returns based on a common factor; the **market return**. For each stock:

$$X_i = \mu_i + w_i Z + \varepsilon_i$$

This is one of the most basic models in finance.

Example: Human Intelligence

Cognitive abilities modeled by latent **intelligence factor**.

This could be further generalized to account for multiple types of intelligence.

Factor Analysis Model

The model assumes the following stochastic representation of $X = (X_1, \dots, X_m)$:

$$X = \mu + WZ + \varepsilon, \quad Z \sim N_r(0, I_r), \quad \varepsilon \sim N_m(0, \Psi), \quad Z \perp\!\!\!\perp \varepsilon,$$

where Ψ is a diagonal covariance matrix.

The latent factors Z

As the two examples suggest, often in this context Z has a specific interpretation.

In PPCA we have the same representation with $\Psi = \sigma^2 I_m$ (isotropic noise).

In FA more emphasis on interpreting the latent factors.

Parametrization and identifiability

$X = \mu + WZ + \varepsilon$ is Gaussian with the induced covariance structure:

$$\Sigma = WW^T + \Psi.$$

→ PPCA: $\Psi = \sigma^2 I_m$

$$(\mu, W, \Psi) \sim (\mu, WU, \Psi) \quad U \in O(m)$$

Lack of identifiability

As for PPCA, W is not uniquely identified.

- ▶ Replacing W with WU for $U \in O(m)$ does not change the distribution.
- ▶ This has important consequences for model interpretability.

Dealing with non-uniqueness of W

Approach 1: Constraint W so that $W^T \Psi^{-1} W$ diagonal.

$$\tilde{X} = \Psi^{-1/2} X$$

- ▶ Multiply $X = \mu + WZ + \varepsilon$ by $\Psi^{-1/2}$ to get $\tilde{X} = \tilde{\mu} + \tilde{W}Z + \tilde{\varepsilon}$ with $\tilde{\varepsilon} \sim N(0, I_m)$.
- ▶ We have $W^T \Psi^{-1} W = \tilde{W}^T \tilde{W}$ so this corresponds to orthogonality of the columns of \tilde{W}^T .

$$\tilde{W} = \Psi^{-1/2} W$$

$$\tilde{\varepsilon} = \Psi^{-1/2} \varepsilon$$

Approach 2: Apply **varimax rotation** for interpretability.

- ▶ Consider any $\widehat{W} \in \mathbb{R}^{m \times r}$. We find U such that $\widehat{W}U$ more interpretable.
- ▶ Define $M \in \mathbb{R}^{m \times r}$ by $M_{ij} = \frac{(WU)_{ij}^2}{\sum_{j=1}^r (WU)_{ij}^2}$ then find the appropriate U **maximizing**:

$$\|M - \frac{1}{m} \mathbf{1}_m \mathbf{1}_m^T M\|_F^2.$$

- ▶ This results with solutions such that each column of M has a bunch of big entries and the remaining ones are negligible.

Fitting the factor analysis model

Data: $\mathbf{x}_1, \dots, \mathbf{x}_n$ from the model $X = \mu + WZ + \varepsilon$, $\varepsilon \sim N(0, \Psi)$.

The most canonical way to estimate the parameters is via the maximum likelihood.

The MLE is not given in a closed form.

- ▶ We could use the EM algorithm.

Alternatively use the fact that MLE has closed form if $\Psi = \sigma^2 I_m$.

Suppose Ψ is known. (fix)

Denote $\tilde{X} = \Psi^{-1/2} X$, $\tilde{\mu} = \Psi^{-1/2} \mu$, $\tilde{W} = \Psi^{-1/2} W$ and $\tilde{\varepsilon} = \Psi^{-1/2} \varepsilon \sim N(0, I_m)$

$$\tilde{X} = \tilde{\mu} + \tilde{W}Z + \tilde{\varepsilon}. \quad (\text{PPCA with } \sigma^2 = 1)$$

Define $\tilde{S}_n = \Psi^{-1/2} S_n \Psi^{-1/2}$ with spectral decomposition $\tilde{S}_n = U \tilde{\Lambda} U^\top$. The MLE:

$$\hat{W} = U_r \Theta R,$$

where:

- ▶ R is any orthogonal matrix, U_r first r columns of U ,
- ▶ Θ is a diagonal matrix with i -th entry equal to $\sqrt{\max\{0, \tilde{\lambda}_i - 1\}}$.

We can now apply this iteratively, where the update on Ψ

Choosing the Number of Factors

$$X = \mu + \cancel{WZ} + \epsilon \quad x_1, \dots, x_n \quad R_n \approx I_m$$

sample correl.

Determining the number of latent factors r is critical in factor analysis.

Overestimating r leads to overfitting, underestimating r leads to loss of structure.

There are several common methods. We focus on Horn's Parallel Analysis (PA).

Key ideas of Horn's Parallel Analysis

If no latent signal, the sample **correlation** matrix should resemble the identity matrix.

Depending on (n, m) the actual eigenvalues can still be far from 1.

PA compares eigenvalues of observed data with those obtained from Monte Carlo simulations of purely random noise.

Horn's Parallel Analysis

$$X = \mu + \varepsilon$$

Based on the observed data $X \in \mathbb{R}^{n \times m}$:

1. Compute eigenvalues $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_m$ of sample correlation matrix R_n , where

$$R_n = D^{-1/2} S_n D^{-1/2}.$$

$$X^{(b)} = \begin{pmatrix} -x_1^{(b)} & \dots & -x_n^{(b)} \\ \vdots & & \vdots \end{pmatrix}$$

2. Generate B simulated datasets $X^{(b)} \in \mathbb{R}^{n \times m}$ from $N_m(\mathbf{0}, I_m)$.
3. Compute sample correlation matrices $R_n^{(b)}$ for each simulated dataset.
4. Compute average null eigenvalues:

$$x_i^{(b)} \sim N(0, I_m)$$

$$\lambda_j^{\text{random}} = \frac{1}{B} \sum_{b=1}^B \lambda_j^{(b)}, \quad j = 1, \dots, m. \quad (2)$$

5. Retain factors where

$$\lambda_j > \lambda_j^{\text{random}}. \quad (3)$$

Application: Factor Analysis on Personality Traits Data

We analyze real survey data from the `bfi` dataset in the `psych` R package.

The dataset consists of 2,800 responses to 25 personality-related questions.

These questions measure the Big Five Personality Traits:

- ▶ Neuroticism (N)
- ▶ Extraversion (E)
- ▶ Conscientiousness (C)
- ▶ Agreeableness (A)
- ▶ Openness (O)

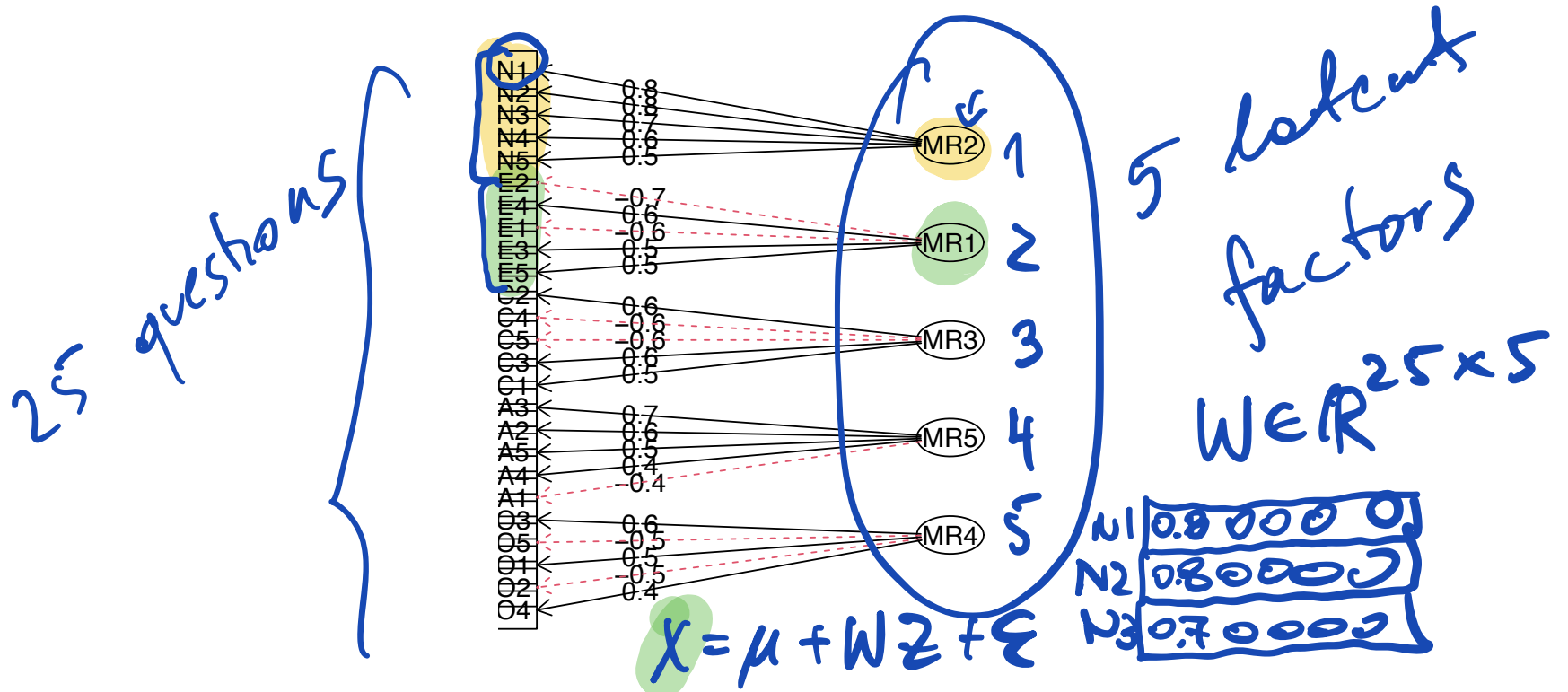
Responses are on the 1-6 scale indicating agreement strength.

We discard demographic variables (gender, age, education).

Factor Analysis Results

PA suggests retaining **five** factors, matching the Big Five personality traits (nice!).

Factor loadings after **varimax rotation** confirm that the extracted factors correspond to the expected latent traits.



Summary

Factor Analysis is a popular method in multivariate statistics.

It is similar to PPCA and it has clear motivating examples.

The lack of identifiability creates a challenge in the interpretation of factor loadings.

Choosing the number of factors (if there is no clear insight) may be also hard.

- ▶ Horn's Parallel Analysis is a simple solution that tends to perform well in practice.

The resulting form of the covariance matrix $\Sigma = WW^T + \Psi$ can be exploited and generalized in many creative ways.