

Applied Statistical Methods II

Some additional topics I

Method of moments (MoM)

In general:

- Let $\mathbf{Y} \in \mathbb{R}^n$ be observed data and $\boldsymbol{\theta} \in \mathbb{R}^p$ be a population parameter of interest
- Define a function $\mathbf{f} : \mathbb{R}^n \times \mathbb{R}^p \rightarrow \mathbb{R}^q$ s.t.

$$E\{\mathbf{f}(\mathbf{Y}, \boldsymbol{\theta})\} = \mathbf{0}_q$$

- Idea: like the score function, the function \mathbf{f} identifies $\boldsymbol{\theta}$.
What is the small value of q that will identify $\boldsymbol{\theta}$?
- Examples:

- Estimating the mean: if $EY_i = \mu$, $f(\mathbf{Y}, \mu) = \sum_{i=1}^n Y_i - n\mu$.

- Variance: if $Y_i \sim (\mu, \sigma^2)$,

$$f(\mathbf{Y}, \sigma^2) = (n-1) - \frac{1}{\sigma^2} \sum_{i=1}^n (Y_i - \bar{Y})^2$$

- Any score function: $\mathbf{f}(\mathbf{Y}, \boldsymbol{\theta}) = \nabla_{\boldsymbol{\theta}} \ell(\boldsymbol{\theta}; \mathbf{Y})$.
- Benefit of MoM: we don't have to know the distribution of \mathbf{Y} ! Only need its moments.

Method of moments (MoM)

In general:

- Let $\mathbf{Y} \in \mathbb{R}^n$ be observed data and $\boldsymbol{\theta} \in \mathbb{R}^p$ be a population parameter of interest
- Define a function $\mathbf{f} : \mathbb{R}^n \times \mathbb{R}^p \rightarrow \mathbb{R}^q$ s.t.

$$E\{\mathbf{f}(\mathbf{Y}, \boldsymbol{\theta})\} = \mathbf{0}_q$$

- Idea: like the score function, the function \mathbf{f} identifies $\boldsymbol{\theta}$.
What is the small value of q that will identify $\boldsymbol{\theta}$?
- Examples:

- Estimating the mean: if $EY_i = \mu$, $f(\mathbf{Y}, \mu) = \sum_{i=1}^n Y_i - n\mu$.

- Variance: if $Y_i \sim (\mu, \sigma^2)$,

$$f(\mathbf{Y}, \sigma^2) = (n-1) - \frac{1}{\sigma^2} \sum_{i=1}^n (Y_i - \bar{Y})^2$$

- Any score function: $\mathbf{f}(\mathbf{Y}, \boldsymbol{\theta}) = \nabla_{\boldsymbol{\theta}} \ell(\boldsymbol{\theta}; \mathbf{Y})$.

- Benefit of MoM: we don't have to know the distribution of \mathbf{Y} ! Only need its moments.

Method of moments (MoM)

In general:

- Let $\mathbf{Y} \in \mathbb{R}^n$ be observed data and $\theta \in \mathbb{R}^p$ be a population parameter of interest
- Define a function $\mathbf{f} : \mathbb{R}^n \times \mathbb{R}^p \rightarrow \mathbb{R}^q$ s.t.

$$E \{ \mathbf{f}(\mathbf{Y}, \theta) \} = \mathbf{0}_q$$

- Idea: like the score function, the function \mathbf{f} identifies θ .
What is the small value of q that will identify θ ?
- Examples:

- Estimating the mean: if $EY_i = \mu$, $f(\mathbf{Y}, \mu) = \sum_{i=1}^n Y_i - n\mu$.

- Variance: if $Y_i \sim (\mu, \sigma^2)$,

$$f(\mathbf{Y}, \sigma^2) = (n-1) - \frac{1}{\sigma^2} \sum_{i=1}^n (Y_i - \bar{Y})^2$$

- Any score function: $\mathbf{f}(\mathbf{Y}, \theta) = \nabla_{\theta} \ell(\theta; \mathbf{Y})$.

- Benefit of MoM: we don't have to know the distribution of \mathbf{Y} ! Only need its moments.

Method of moments (MoM)

In general:

- Let $\mathbf{Y} \in \mathbb{R}^n$ be observed data and $\boldsymbol{\theta} \in \mathbb{R}^p$ be a population parameter of interest
- Define a function $\mathbf{f} : \mathbb{R}^n \times \mathbb{R}^p \rightarrow \mathbb{R}^q$ s.t.

$$E \{ \mathbf{f}(\mathbf{Y}, \boldsymbol{\theta}) \} = \mathbf{0}_q$$

- Idea: like the score function, the function \mathbf{f} identifies $\boldsymbol{\theta}$.
What is the small value of q that will identify $\boldsymbol{\theta}$?
- Examples:

- Estimating the mean: if $EY_i = \mu$, $f(\mathbf{Y}, \mu) = \sum_{i=1}^n Y_i - n\mu$.

- Variance: if $Y_i \sim (\mu, \sigma^2)$,
 $f(\mathbf{Y}, \sigma^2) = (n-1) - \frac{1}{\sigma^2} \sum_{i=1}^n (Y_i - \bar{Y})^2$

- Any score function: $\mathbf{f}(\mathbf{Y}, \boldsymbol{\theta}) = \nabla_{\boldsymbol{\theta}} \ell(\boldsymbol{\theta}; \mathbf{Y})$.

- Benefit of MoM: we don't have to know the distribution of \mathbf{Y} ! Only need its moments.

Method of moments (MoM)

In general:

- Let $\mathbf{Y} \in \mathbb{R}^n$ be observed data and $\boldsymbol{\theta} \in \mathbb{R}^p$ be a population parameter of interest
- Define a function $\mathbf{f} : \mathbb{R}^n \times \mathbb{R}^p \rightarrow \mathbb{R}^q$ s.t.

$$E \{ \mathbf{f}(\mathbf{Y}, \boldsymbol{\theta}) \} = \mathbf{0}_q$$

- Idea: like the score function, the function \mathbf{f} identifies $\boldsymbol{\theta}$.
What is the small value of q that will identify $\boldsymbol{\theta}$?
- Examples:

- Estimating the mean: if $EY_i = \mu$, $f(\mathbf{Y}, \mu) = \sum_{i=1}^n Y_i - n\mu$.

- Variance: if $Y_i \sim (\mu, \sigma^2)$,

$$f(\mathbf{Y}, \sigma^2) = (n-1) - \frac{1}{\sigma^2} \sum_{i=1}^n (Y_i - \bar{Y})^2$$

- Any score function: $\mathbf{f}(\mathbf{Y}, \boldsymbol{\theta}) = \nabla_{\boldsymbol{\theta}} \ell(\boldsymbol{\theta}; \mathbf{Y})$.
- Benefit of MoM: we don't have to know the distribution of \mathbf{Y} ! Only need its moments.

MINQUE (C.R. Rao, 1973)

Suppose $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon} \in \mathbb{R}^n$, $\boldsymbol{\epsilon} \sim (0, \sum_{r=1}^b \theta_r \mathbf{B}_r)$.

- Estimate θ_r with method of moments.
- We'll use MINQUE: Minimum Norm Quadratic Unbiased Estimation
- Idea: For $\mathbf{A}_s \in \mathbb{R}^{n \times n}$,

$$E(\mathbf{Y}^T \mathbf{A}_s \mathbf{Y}) = \boldsymbol{\beta}^T \mathbf{X}^T \mathbf{A}_s \mathbf{X} \boldsymbol{\beta} + \sum_{r=1}^b \theta_r \text{Tr}(\mathbf{B}_r \mathbf{A}_s)$$

- We will set $\hat{\theta}_s = \mathbf{Y}^T \mathbf{A}_s \mathbf{Y}$. What conditions do we need on \mathbf{A}_s such that $E(\hat{\theta}_s) = \theta_s$?
- $\mathbf{X}^T \mathbf{A}_s \mathbf{X} = \mathbf{0}_p$.
- $\text{Tr}(\mathbf{B}_r \mathbf{A}_s) = 1 \{r = s\}$.
- Is \mathbf{A}_s symmetric? How about positive semi-definite?
- Question: how do we choose $\mathbf{A}_1, \dots, \mathbf{A}_b$?

MINQUE (C.R. Rao, 1973)

Suppose $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon} \in \mathbb{R}^n$, $\boldsymbol{\epsilon} \sim (0, \sum_{r=1}^b \theta_r \mathbf{B}_r)$.

- Estimate θ_r with method of moments.
- We'll use MINQUE: Minimum Norm Quadratic Unbiased Estimation
- Idea: For $\mathbf{A}_s \in \mathbb{R}^{n \times n}$,

$$E(\mathbf{Y}^T \mathbf{A}_s \mathbf{Y}) = \boldsymbol{\beta}^T \mathbf{X}^T \mathbf{A}_s \mathbf{X} \boldsymbol{\beta} + \sum_{r=1}^b \theta_r \text{Tr}(\mathbf{B}_r \mathbf{A}_s)$$

- We will set $\hat{\theta}_s = \mathbf{Y}^T \mathbf{A}_s \mathbf{Y}$. What conditions do we need on \mathbf{A}_s such that $E(\hat{\theta}_s) = \theta_s$?
- $\mathbf{X}^T \mathbf{A}_s \mathbf{X} = \mathbf{0}_p$.
- $\text{Tr}(\mathbf{B}_r \mathbf{A}_s) = 1 \{r = s\}$.
- Is \mathbf{A}_s symmetric? How about positive semi-definite?
- Question: how do we choose $\mathbf{A}_1, \dots, \mathbf{A}_b$?

MINQUE (C.R. Rao, 1973)

Suppose $\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon} \in \mathbb{R}^n$, $\boldsymbol{\epsilon} \sim (0, \sum_{r=1}^b \theta_r \mathbf{B}_r)$.

- Estimate θ_r with method of moments.
- We'll use MINQUE: Minimum Norm Quadratic Unbiased Estimation
- Idea: For $\mathbf{A}_s \in \mathbb{R}^{n \times n}$,

$$E(\mathbf{Y}^T \mathbf{A}_s \mathbf{Y}) = \boldsymbol{\beta}^T \mathbf{X}^T \mathbf{A}_s \mathbf{X} \boldsymbol{\beta} + \sum_{r=1}^b \theta_r \text{Tr}(\mathbf{B}_r \mathbf{A}_s)$$

- We will set $\hat{\theta}_s = \mathbf{Y}^T \mathbf{A}_s \mathbf{Y}$. What conditions do we need on \mathbf{A}_s such that $E(\hat{\theta}_s) = \theta_s$?
- $\mathbf{X}^T \mathbf{A}_s \mathbf{X} = \mathbf{0}_p$.
- $\text{Tr}(\mathbf{B}_r \mathbf{A}_s) = 1 \{r = s\}$.
- Is \mathbf{A}_s symmetric? How about positive semi-definite?
- Question: how do we choose $\mathbf{A}_1, \dots, \mathbf{A}_b$?

MINQUE (cont.)

$$\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon} \in \mathbb{R}^n, \boldsymbol{\epsilon} \sim (0, \sum_{r=1}^b \theta_r \mathbf{B}_r)$$

- $\hat{\theta}_s = \mathbf{Y}^T \mathbf{A}_s \mathbf{Y}$, $\mathbf{X}^T \mathbf{A}_s \mathbf{X} = \mathbf{0}_p$, $\text{Tr}(\mathbf{B}_r \mathbf{A}_s) = 1 \{r = s\}$
- We'll choose \mathbf{A}_s such that $\hat{\theta}_s$ has minimal variance.
- Under assumptions of normality:
$$\text{Var}(\hat{\theta}_s) = 2 \sum_{r,t} \theta_r \theta_t \text{Tr}(\mathbf{A}_s \mathbf{B}_r \mathbf{A}_s \mathbf{B}_t)$$
- Problem: $\text{Var}(\hat{\theta}_s)$ depends on $\boldsymbol{\theta}$, the parameter we're trying to estimate!
- Any ideas how to circumvent this?
 - 1 Get a consistent, but inefficient, estimator $\hat{\boldsymbol{\theta}}^{(0)}$ using $\mathbf{A}_1^{(0)}, \dots, \mathbf{A}_b^{(0)}$.
 - 2 Re-compute $\mathbf{A}_1, \dots, \mathbf{A}_r$ by minimizing $\text{Var}(\hat{\theta}_1^{(0)}), \dots, \text{Var}(\hat{\theta}_b^{(0)})$.

MINQUE (cont.)

$$\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon} \in \mathbb{R}^n, \boldsymbol{\epsilon} \sim (0, \sum_{r=1}^b \theta_r \mathbf{B}_r)$$

- $\hat{\theta}_s = \mathbf{Y}^T \mathbf{A}_s \mathbf{Y}$, $\mathbf{X}^T \mathbf{A}_s \mathbf{X} = \mathbf{0}_p$, $\text{Tr}(\mathbf{B}_r \mathbf{A}_s) = 1 \{r = s\}$
- We'll choose \mathbf{A}_s such that $\hat{\theta}_s$ has minimal variance.
- Under assumptions of normality:
$$\text{Var}(\hat{\theta}_s) = 2 \sum_{r,t} \theta_r \theta_t \text{Tr}(\mathbf{A}_s \mathbf{B}_r \mathbf{A}_s \mathbf{B}_t)$$
- Problem: $\text{Var}(\hat{\theta}_s)$ depends on $\boldsymbol{\theta}$, the parameter we're trying to estimate!
- Any ideas how to circumvent this?
 - 1 Get a consistent, but inefficient, estimator $\hat{\boldsymbol{\theta}}^{(0)}$ using $\mathbf{A}_1^{(0)}, \dots, \mathbf{A}_b^{(0)}$.
 - 2 Re-compute $\mathbf{A}_1, \dots, \mathbf{A}_r$ by minimizing $\text{Var}(\hat{\theta}_1^{(0)}), \dots, \text{Var}(\hat{\theta}_b^{(0)})$.

MINQUE (cont.)

$$\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon} \in \mathbb{R}^n, \boldsymbol{\epsilon} \sim (0, \sum_{r=1}^b \theta_r \mathbf{B}_r)$$

- $\hat{\theta}_s = \mathbf{Y}^T \mathbf{A}_s \mathbf{Y}$, $\mathbf{X}^T \mathbf{A}_s \mathbf{X} = \mathbf{0}_p$, $\text{Tr}(\mathbf{B}_r \mathbf{A}_s) = 1 \{r = s\}$
- We'll choose \mathbf{A}_s such that $\hat{\theta}_s$ has minimal variance.
- Under assumptions of normality:
$$\text{Var}(\hat{\theta}_s) = 2 \sum_{r,t} \theta_r \theta_t \text{Tr}(\mathbf{A}_s \mathbf{B}_r \mathbf{A}_s \mathbf{B}_t)$$
- Problem: $\text{Var}(\hat{\theta}_s)$ depends on $\boldsymbol{\theta}$, the parameter we're trying to estimate!
- Any ideas how to circumvent this?
 - 1 Get a consistent, but inefficient, estimator $\hat{\boldsymbol{\theta}}^{(0)}$ using $\mathbf{A}_1^{(0)}, \dots, \mathbf{A}_b^{(0)}$.
 - 2 Re-compute $\mathbf{A}_1, \dots, \mathbf{A}_r$ by minimizing $\text{Var}(\hat{\theta}_1^{(0)}), \dots, \text{Var}(\hat{\theta}_b^{(0)})$.

Generalized method of moments

- Let $\mathbf{Y} \in \mathbb{R}^n$ be observed data and $\theta \in \mathbb{R}^p$ be a population parameter of interest.
- Specify the function $\mathbf{f} : \mathbb{R}^n \times \mathbb{R}^p \rightarrow \mathbb{R}^q$ such that $E\{\mathbf{f}(\mathbf{Y}, \theta)\} = \mathbf{0}_q$.
- Question: how do we estimate and perform inference on θ when the distribution of \mathbf{Y} is unknown?
- The problem was solved for fixed dimensions p and q in Hansen (1982), which won him the Nobel Prize in economics in 2013.
- Let $\mathbf{W} \in \mathbb{R}^{q \times q}$ be p.d. By moment condition, we consider the class of estimators

$$\hat{\theta}_W = \underset{\theta}{\operatorname{argmin}} \mathbf{f}(\mathbf{Y}, \theta)^T \mathbf{W} \mathbf{f}(\mathbf{Y}, \theta).$$

- We require $q \geq p$ (Why?) If $q = p$, $\mathbf{f}(\mathbf{Y}, \hat{\theta}_W) = \mathbf{0}_p$, and $\hat{\theta}_W$ is invariant to \mathbf{W} .

Generalized method of moments

- Let $\mathbf{Y} \in \mathbb{R}^n$ be observed data and $\theta \in \mathbb{R}^p$ be a population parameter of interest.
- Specify the function $\mathbf{f} : \mathbb{R}^n \times \mathbb{R}^p \rightarrow \mathbb{R}^q$ such that $E\{\mathbf{f}(\mathbf{Y}, \theta)\} = \mathbf{0}_q$.
- Question: how do we estimate and perform inference on θ when the distribution of \mathbf{Y} is unknown?
- The problem was solved for fixed dimensions p and q in Hansen (1982), which won him the Nobel Prize in economics in 2013.
- Let $\mathbf{W} \in \mathbb{R}^{q \times q}$ be p.d. By moment condition, we consider the class of estimators

$$\hat{\theta}_W = \underset{\theta}{\operatorname{argmin}} \mathbf{f}(\mathbf{Y}, \theta)^T \mathbf{W} \mathbf{f}(\mathbf{Y}, \theta).$$

- We require $q \geq p$ (Why?) If $q = p$, $\mathbf{f}(\mathbf{Y}, \hat{\theta}_W) = \mathbf{0}_p$, and $\hat{\theta}_W$ is invariant to \mathbf{W} .

Generalized method of moments

$$E\{\mathbf{f}(\mathbf{Y}, \theta)\} = \mathbf{0}_q,$$

$$\hat{\theta}_W = \operatorname{argmin}_{\theta} \mathbf{f}(\mathbf{Y}, \theta)^T \mathbf{W} \mathbf{f}(\mathbf{Y}, \theta)$$

- If $q > p$, then $\mathbf{f}(\mathbf{Y}, \theta) \neq 0$ for any θ . How do we select \mathbf{W} ?
- By same technique as MINQUE: Minimize the asymptotic variance of $\hat{\theta}_W$ (which depends on θ).
 - 1 Start with $\mathbf{W} = \mathbf{I}_q$. $\hat{\theta}_{\mathbf{I}_q}$ is a \sqrt{n} -consistent estimator for θ .
 - 2 Set $\mathbf{W} = \mathbf{W}(\hat{\theta}_{\mathbf{I}_q})$ to be the optimal \mathbf{W} .
 - 3 Using this \mathbf{W} , $\hat{\theta}_W$ has the smallest possible asymptotic variance!
- Called **Hansen's two-step estimator**. Used all the time in economics, and recently to study non-random missing data in mass spectrometry data (McKenna et al., 2020).

Generalized method of moments

$$E\{\mathbf{f}(\mathbf{Y}, \theta)\} = \mathbf{0}_q,$$

$$\hat{\theta}_W = \operatorname{argmin}_{\theta} \mathbf{f}(\mathbf{Y}, \theta)^T \mathbf{W} \mathbf{f}(\mathbf{Y}, \theta)$$

- If $q > p$, then $\mathbf{f}(\mathbf{Y}, \theta) \neq 0$ for any θ . How do we select \mathbf{W} ?
- By same technique as MINQUE: Minimize the asymptotic variance of $\hat{\theta}_W$ (which depends on θ).
 - 1 Start with $\mathbf{W} = \mathbf{I}_q$. $\hat{\theta}_{\mathbf{I}_q}$ is a \sqrt{n} -consistent estimator for θ .
 - 2 Set $\mathbf{W} = \mathbf{W}(\hat{\theta}_{\mathbf{I}_q})$ to be the optimal \mathbf{W} .
 - 3 Using this \mathbf{W} , $\hat{\theta}_W$ has the smallest possible asymptotic variance!
- Called **Hansen's two-step estimator**. Used all the time in economics, and recently to study non-random missing data in mass spectrometry data (McKenna et al., 2020).

Generalized method of moments

$$E\{\mathbf{f}(\mathbf{Y}, \theta)\} = \mathbf{0}_q,$$

$$\hat{\theta}_W = \operatorname{argmin}_{\theta} \mathbf{f}(\mathbf{Y}, \theta)^T \mathbf{W} \mathbf{f}(\mathbf{Y}, \theta)$$

- If $q > p$, then $\mathbf{f}(\mathbf{Y}, \theta) \neq 0$ for any θ . How do we select \mathbf{W} ?
- By same technique as MINQUE: Minimize the asymptotic variance of $\hat{\theta}_W$ (which depends on θ).
 - 1 Start with $\mathbf{W} = \mathbf{I}_q$. $\hat{\theta}_{\mathbf{I}_q}$ is a \sqrt{n} -consistent estimator for θ .
 - 2 Set $\mathbf{W} = \mathbf{W}(\hat{\theta}_{\mathbf{I}_q})$ to be the optimal \mathbf{W} .
 - 3 Using this \mathbf{W} , $\hat{\theta}_W$ has the smallest possible asymptotic variance!
- Called **Hansen's two-step estimator**. Used all the time in economics, and recently to study non-random missing data in mass spectrometry data (McKenna et al., 2020).

Factor Analysis and dimension reduction in noisy data

Overview

- Likely the most important tool in all of high dimensional statistics.
- Used in every discipline: machine learning, genetics, economics, sociology, psychology, education, political science, and many others...
- Comes in many flavors: linear, non-linear, finite dimensional, infinite dimensional (i.e. factor analysis in a RKHS)
- Core idea: find a small number of **factors** that explain most of the variation in the data.
- In the machine learning/computing science community: often treat this as an optimization problem.
 - Problem: we have no way of understanding the statistical uncertainty in our estimates.
- We will study this in the context of a statistical model.
 - This will allow us to justify our choice of estimators.
 - Study the statistical uncertainty.

Overview

- Likely the most important tool in all of high dimensional statistics.
- Used in every discipline: machine learning, genetics, economics, sociology, psychology, education, political science, and many others...
- Comes in many flavors: linear, non-linear, finite dimensional, infinite dimensional (i.e. factor analysis in a RKHS)
- Core idea: find a small number of **factors** that explain most of the variation in the data.
- In the machine learning/computing science community: often treat this as an optimization problem.
 - Problem: we have no way of understanding the statistical uncertainty in our estimates.
- We will study this in the context of a statistical model.
 - This will allow us to justify our choice of estimators.
 - Study the statistical uncertainty.

Linear factor analysis in noisy data I

- $\mathbf{y}_g = \mathbf{C}\ell_g + \mathbf{e}_g \in \mathbb{R}^n$ is the expression of gene $g = 1, \dots, p$ (or brain region g , survey question g , etc) across n samples.
 - Will assume that $\mathbf{e}_g \sim (\mathbf{0}, \sigma_g^2 I_n)$ and \mathbf{e}_g is independent of \mathbf{e}_h
 - p and n are large. Think $p \gtrsim n$, and maybe $p \gg n$.
- The **factors** $\mathbf{C} \in \mathbb{R}^{n \times K}$ are **shared** across all genes $g = 1, \dots, p$.
- The **loadings** $\ell_g \in \mathbb{R}^K$ are the effects of \mathbf{C} on the expression of gene g . Examples:
 - Some of the columns of \mathbf{C} might correspond to biological factors like cell type. If a person has more T cells, their expression will look different from someone with more B cells.
 - Maybe some are related to disease (sick vs. healthy), technical factors (processing batch), etc.

Linear factor analysis in noisy data II

- $K \ll \min(n, p)$, i.e. want to explain the variation in $\mathbf{y}_1, \dots, \mathbf{y}_p$ with a small number of factors.
- \mathbf{C} can induce dependencies across genes. Assuming rows $i = 1, \dots, n$ of \mathbf{C} are i.i.d with $\text{Var}(\mathbf{C}_i) = \Psi$:

$$\text{Cov}(y_{gi}, y_{hi}) = \ell_g^T \Psi \ell_h, \quad g \neq h = 1, \dots, p$$

- We only observe \mathbf{y}_g . Our goal is to recover ℓ_1, \dots, ℓ_p and \mathbf{C} .

Fundamental problem in factor analysis

$$\mathbf{y}_g = \mathbf{C}\ell_g + \mathbf{e}_g \in \mathbb{R}^n \text{ for } g = 1, \dots, p$$

$$\mathbf{Y}_{p \times n} = \begin{bmatrix} \mathbf{y}_1^T \\ \vdots \\ \mathbf{y}_p^T \end{bmatrix} = \underbrace{\begin{pmatrix} \ell_1^T \\ \vdots \\ \ell_p^T \end{pmatrix}}_{\mathbf{L}_{p \times K}} \underbrace{\begin{pmatrix} \mathbf{c}_1 & \cdots & \mathbf{c}_n \end{pmatrix}}_{\mathbf{C}^T} + \underbrace{\begin{pmatrix} \mathbf{e}_1^T \\ \vdots \\ \mathbf{e}_p^T \end{pmatrix}}_{\mathbf{E}_{p \times n}}$$

- Recall: $p \gtrsim n$, maybe $p \gg n$.
- Recall: we only observe \mathbf{Y} . Goal is to recover $\mathbf{C} \in \mathbb{R}^{n \times K}$ and $\mathbf{L} \in \mathbb{R}^{p \times K}$, $K \ll \min(n, p)$
- Problem: without further assumptions, \mathbf{L} and \mathbf{C} are not identifiable! Why? We can only recover \mathbf{LC}^T :

$$E(\mathbf{Y} | \mathbf{C}) = \mathbf{LC}^T = \mathbf{LC}^T = \mathbf{LR}(\mathbf{CR}^{-T})^T, \quad \forall \text{ invertible } \mathbf{R} \in \mathbb{R}^{K \times K}$$

- Are $\text{im}(\mathbf{C})$ or $\text{im}(\mathbf{L})$ identifiable?
- Different types of factor analyses place different assumptions on \mathbf{L} and \mathbf{C} to make them identifiable/interpretable:

Fundamental problem in factor analysis

$$\mathbf{y}_g = \mathbf{C}\ell_g + \mathbf{e}_g \in \mathbb{R}^n \text{ for } g = 1, \dots, p$$

$$\mathbf{Y}_{p \times n} = \begin{bmatrix} \mathbf{y}_1^T \\ \vdots \\ \mathbf{y}_p^T \end{bmatrix} = \underbrace{\begin{pmatrix} \ell_1^T \\ \vdots \\ \ell_p^T \end{pmatrix}}_{\mathbf{L}_{p \times K}} \underbrace{\begin{pmatrix} \mathbf{c}_1 & \cdots & \mathbf{c}_n \end{pmatrix}}_{\mathbf{C}^T} + \underbrace{\begin{pmatrix} \mathbf{e}_1^T \\ \vdots \\ \mathbf{e}_p^T \end{pmatrix}}_{\mathbf{E}_{p \times n}}$$

- Recall: $p \gtrsim n$, maybe $p \gg n$.
- Recall: we only observe \mathbf{Y} . Goal is to recover $\mathbf{C} \in \mathbb{R}^{n \times K}$ and $\mathbf{L} \in \mathbb{R}^{p \times K}$, $K \ll \min(n, p)$
- Problem: without further assumptions, \mathbf{L} and \mathbf{C} are not identifiable! Why? We can only recover \mathbf{LC}^T :

$$E(\mathbf{Y} | \mathbf{C}) = \mathbf{LC}^T = \mathbf{LC}^T = \mathbf{LR}(\mathbf{CR}^{-T})^T, \quad \forall \text{ invertible } \mathbf{R} \in \mathbb{R}^{K \times K}$$

- Are $\text{im}(\mathbf{C})$ or $\text{im}(\mathbf{L})$ identifiable?
- Different types of factor analyses place different assumptions on \mathbf{L} and \mathbf{C} to make them identifiable/interpretable:

Fundamental problem in factor analysis

$$\mathbf{y}_g = \mathbf{C}\ell_g + \mathbf{e}_g \in \mathbb{R}^n \text{ for } g = 1, \dots, p$$

$$\mathbf{Y}_{p \times n} = \begin{bmatrix} \mathbf{y}_1^T \\ \vdots \\ \mathbf{y}_p^T \end{bmatrix} = \underbrace{\begin{pmatrix} \ell_1^T \\ \vdots \\ \ell_p^T \end{pmatrix}}_{\mathbf{L}_{p \times K}} \underbrace{\begin{pmatrix} \mathbf{c}_1 & \cdots & \mathbf{c}_n \end{pmatrix}}_{\mathbf{C}^T} + \underbrace{\begin{pmatrix} \mathbf{e}_1^T \\ \vdots \\ \mathbf{e}_p^T \end{pmatrix}}_{\mathbf{E}_{p \times n}}$$

- Recall: $p \gtrsim n$, maybe $p \gg n$.
- Recall: we only observe \mathbf{Y} . Goal is to recover $\mathbf{C} \in \mathbb{R}^{n \times K}$ and $\mathbf{L} \in \mathbb{R}^{p \times K}$, $K \ll \min(n, p)$
- Problem: without further assumptions, \mathbf{L} and \mathbf{C} are not identifiable! Why? We can only recover \mathbf{LC}^T :

$$E(\mathbf{Y} | \mathbf{C}) = \mathbf{LC}^T = \mathbf{LC}^T = \mathbf{LR}(\mathbf{CR}^{-T})^T, \quad \forall \text{ invertible } \mathbf{R} \in \mathbb{R}^{K \times K}$$

- Are $\text{im}(\mathbf{C})$ or $\text{im}(\mathbf{L})$ identifiable?
- Different types of factor analyses place different assumptions on \mathbf{L} and \mathbf{C} to make them identifiable/interpretable:

Fundamental problem in factor analysis

$$\mathbf{y}_g = \mathbf{C}\ell_g + \mathbf{e}_g \in \mathbb{R}^n \text{ for } g = 1, \dots, p$$

$$\mathbf{Y}_{p \times n} = \begin{bmatrix} \mathbf{y}_1^T \\ \vdots \\ \mathbf{y}_p^T \end{bmatrix} = \underbrace{\begin{pmatrix} \ell_1^T \\ \vdots \\ \ell_p^T \end{pmatrix}}_{\mathbf{L}_{p \times K}} \underbrace{\begin{pmatrix} \mathbf{c}_1 & \cdots & \mathbf{c}_n \end{pmatrix}}_{\mathbf{C}^T} + \underbrace{\begin{pmatrix} \mathbf{e}_1^T \\ \vdots \\ \mathbf{e}_p^T \end{pmatrix}}_{\mathbf{E}_{p \times n}}$$

- Recall: $p \gtrsim n$, maybe $p \gg n$.
- Recall: we only observe \mathbf{Y} . Goal is to recover $\mathbf{C} \in \mathbb{R}^{n \times K}$ and $\mathbf{L} \in \mathbb{R}^{p \times K}$, $K \ll \min(n, p)$
- Problem: without further assumptions, \mathbf{L} and \mathbf{C} are not identifiable! Why? We can only recover \mathbf{LC}^T :

$$E(\mathbf{Y} | \mathbf{C}) = \mathbf{LC}^T = \mathbf{LC}^T = \mathbf{LR}(\mathbf{CR}^{-T})^T, \quad \forall \text{ invertible } \mathbf{R} \in \mathbb{R}^{K \times K}$$

- Are $\text{im}(\mathbf{C})$ or $\text{im}(\mathbf{L})$ identifiable?
- Different types of factor analyses place different assumptions on \mathbf{L} and \mathbf{C} to make them identifiable/interpretable:

Fundamental problem in factor analysis

$$\mathbf{y}_g = \mathbf{C}\ell_g + \mathbf{e}_g \in \mathbb{R}^n \text{ for } g = 1, \dots, p$$

$$\mathbf{Y}_{p \times n} = \begin{bmatrix} \mathbf{y}_1^T \\ \vdots \\ \mathbf{y}_p^T \end{bmatrix} = \underbrace{\begin{pmatrix} \ell_1^T \\ \vdots \\ \ell_p^T \end{pmatrix}}_{\mathbf{L}_{p \times K}} \underbrace{\begin{pmatrix} \mathbf{c}_1 & \cdots & \mathbf{c}_n \end{pmatrix}}_{\mathbf{C}^T} + \underbrace{\begin{pmatrix} \mathbf{e}_1^T \\ \vdots \\ \mathbf{e}_p^T \end{pmatrix}}_{\mathbf{E}_{p \times n}}$$

- Recall: $p \gtrsim n$, maybe $p \gg n$.
- Recall: we only observe \mathbf{Y} . Goal is to recover $\mathbf{C} \in \mathbb{R}^{n \times K}$ and $\mathbf{L} \in \mathbb{R}^{p \times K}$, $K \ll \min(n, p)$
- Problem: without further assumptions, \mathbf{L} and \mathbf{C} are not identifiable! Why? We can only recover \mathbf{LC}^T :

$$E(\mathbf{Y} | \mathbf{C}) = \mathbf{LC}^T = \mathbf{LC}^T = \mathbf{LR}(\mathbf{CR}^{-T})^T, \quad \forall \text{ invertible } \mathbf{R} \in \mathbb{R}^{K \times K}$$

- Are $\text{im}(\mathbf{C})$ or $\text{im}(\mathbf{L})$ identifiable?
- Different types of factor analyses place different assumptions on \mathbf{L} and \mathbf{C} to make them identifiable/interpretable:

Fundamental problem in factor analysis

$$\mathbf{y}_g = \mathbf{C}\ell_g + \mathbf{e}_g \in \mathbb{R}^n \text{ for } g = 1, \dots, p$$

$$\mathbf{Y}_{p \times n} = \begin{bmatrix} \mathbf{y}_1^T \\ \vdots \\ \mathbf{y}_p^T \end{bmatrix} = \underbrace{\begin{pmatrix} \ell_1^T \\ \vdots \\ \ell_p^T \end{pmatrix}}_{\mathbf{L}_{p \times K}} \underbrace{\begin{pmatrix} \mathbf{c}_1 & \cdots & \mathbf{c}_n \end{pmatrix}}_{\mathbf{C}^T} + \underbrace{\begin{pmatrix} \mathbf{e}_1^T \\ \vdots \\ \mathbf{e}_p^T \end{pmatrix}}_{\mathbf{E}_{p \times n}}$$

- Recall: $p \gtrsim n$, maybe $p \gg n$.
- Recall: we only observe \mathbf{Y} . Goal is to recover $\mathbf{C} \in \mathbb{R}^{n \times K}$ and $\mathbf{L} \in \mathbb{R}^{p \times K}$, $K \ll \min(n, p)$
- Problem: without further assumptions, \mathbf{L} and \mathbf{C} are not identifiable! Why? We can only recover \mathbf{LC}^T :

$$E(\mathbf{Y} | \mathbf{C}) = \mathbf{LC}^T = \mathbf{LC}^T = \mathbf{LR}(\mathbf{CR}^{-T})^T, \quad \forall \text{ invertible } \mathbf{R} \in \mathbb{R}^{K \times K}$$

- Are $\text{im}(\mathbf{C})$ or $\text{im}(\mathbf{L})$ identifiable?
- Different types of factor analyses place different assumptions on \mathbf{L} and \mathbf{C} to make them identifiable/interpretable.

Fundamental problem in factor analysis

$$\mathbf{y}_g = \mathbf{C}\ell_g + \mathbf{e}_g \in \mathbb{R}^n \text{ for } g = 1, \dots, p$$

$$\mathbf{Y}_{p \times n} = \begin{bmatrix} \mathbf{y}_1^T \\ \vdots \\ \mathbf{y}_p^T \end{bmatrix} = \underbrace{\begin{pmatrix} \ell_1^T \\ \vdots \\ \ell_p^T \end{pmatrix}}_{\mathbf{L}_{p \times K}} \underbrace{\begin{pmatrix} \mathbf{c}_1 & \cdots & \mathbf{c}_n \end{pmatrix}}_{\mathbf{C}^T} + \underbrace{\begin{pmatrix} \mathbf{e}_1^T \\ \vdots \\ \mathbf{e}_p^T \end{pmatrix}}_{\mathbf{E}_{p \times n}}$$

- Recall: $p \gtrsim n$, maybe $p \gg n$.
- Recall: we only observe \mathbf{Y} . Goal is to recover $\mathbf{C} \in \mathbb{R}^{n \times K}$ and $\mathbf{L} \in \mathbb{R}^{p \times K}$, $K \ll \min(n, p)$
- Problem: without further assumptions, \mathbf{L} and \mathbf{C} are not identifiable! Why? We can only recover \mathbf{LC}^T :

$$E(\mathbf{Y} | \mathbf{C}) = \mathbf{LC}^T = \mathbf{LC}^T = \mathbf{LR}(\mathbf{CR}^{-T})^T, \quad \forall \text{ invertible } \mathbf{R} \in \mathbb{R}^{K \times K}$$

- Are $\text{im}(\mathbf{C})$ or $\text{im}(\mathbf{L})$ identifiable?
- Different types of factor analyses place different assumptions on \mathbf{L} and \mathbf{C} to make them identifiable/interpretable.

PCA in noisy data I

$\mathbf{Y}_{p \times n} = \mathbf{L}_{p \times K} \mathbf{C}_{n \times K}^T + \mathbf{E}_{p \times n}$, entries of \mathbf{E} are independent,
 $E_{gi} \sim (0, \sigma_g^2)$.

- PCA can be interpreted as a particular parametrization of \mathbf{L} and \mathbf{C}
- Recall goal of PCA: identify most important sources of variation of \mathbf{Y}
 - PCs can be ordered by corresponding eigenvalue.

$\mathbf{Y}_{p \times n} = \mathbf{L}_{p \times K} \mathbf{C}_{n \times K}^T + \mathbf{E}_{p \times n}$, entries of \mathbf{E} are independent,
 $E_{gi} \sim (0, \sigma_g^2)$.

- PCA can be interpreted as a particular parametrization of \mathbf{L} and \mathbf{C}
- Recall goal of PCA: identify most important sources of variation of \mathbf{Y}
 - PCs can be ordered by corresponding eigenvalue.

PCA in noisy data II

$\mathbf{Y}_{p \times n} = \mathbf{L}_{p \times K} \mathbf{C}_{n \times K}^T + \mathbf{E}_{p \times n}$, entries of \mathbf{E} are independent,
 $E_{gi} \sim (0, \sigma_g^2)$

- Claim: there exists a parametrization of \mathbf{L} and \mathbf{C} s.t.

① $E(\mathbf{Y} | \mathbf{C}) = \mathbf{L} \mathbf{C}^T$.

② $n^{-1} \mathbf{C}^T \mathbf{C} = \mathbf{I}_K$ and $p^{-1} \mathbf{L}^T \mathbf{L} = \text{diag}(\lambda_1, \dots, \lambda_K)$,
 $\lambda_1 \geq \dots \geq \lambda_K > 0$.

- To identify $\mathbf{C}_{\cdot 1}, \dots, \mathbf{C}_{\cdot K}$ up to sign, need to assume
 $\lambda_1 > \dots > \lambda_K > 0$. Why?

③ λ_k and $n^{-1/2} \mathbf{C}_{\cdot k}$ are the k th eigenvalue and eigenvector of
 $n^{-1} \mathbf{C} (p^{-1} \mathbf{L}^T \mathbf{L}) \mathbf{C}^T$

- To show this: compute the singular value decomposition of $E(\mathbf{Y} | \mathbf{C}) = \mathbf{L} \mathbf{C}^T$!
- k th factor is $\mathbf{C}_{\cdot k}$; has loading $\mathbf{L}_{\cdot k}$ with $p^{-1} \mathbf{L}_{\cdot k}^T \mathbf{L}_{\cdot k} = \lambda_k$. Here λ_k is the average effect size for the k th factor.
- Interpretation: $\mathbf{C}_{\cdot k}$ is the k th (out of K) most important factor.
- Goal: estimate $\mathbf{C}_{\cdot k}$ and $\mathbf{L}_{\cdot k}$.

PCA in noisy data II

$\mathbf{Y}_{p \times n} = \mathbf{L}_{p \times K} \mathbf{C}_{n \times K}^T + \mathbf{E}_{p \times n}$, entries of \mathbf{E} are independent,
 $E_{gi} \sim (0, \sigma_g^2)$

- Claim: there exists a parametrization of \mathbf{L} and \mathbf{C} s.t.

① $E(\mathbf{Y} | \mathbf{C}) = \mathbf{L} \mathbf{C}^T$.

② $n^{-1} \mathbf{C}^T \mathbf{C} = \mathbf{I}_K$ and $p^{-1} \mathbf{L}^T \mathbf{L} = \text{diag}(\lambda_1, \dots, \lambda_K)$,
 $\lambda_1 \geq \dots \geq \lambda_K > 0$.

- To identify $\mathbf{C}_{\cdot 1}, \dots, \mathbf{C}_{\cdot K}$ up to sign, need to assume
 $\lambda_1 > \dots > \lambda_K > 0$. Why?

③ λ_k and $n^{-1/2} \mathbf{C}_{\cdot k}$ are the k th eigenvalue and eigenvector of
 $n^{-1} \mathbf{C} (p^{-1} \mathbf{L}^T \mathbf{L}) \mathbf{C}^T$

- To show this: compute the singular value decomposition of $E(\mathbf{Y} | \mathbf{C}) = \mathbf{L} \mathbf{C}^T$!
- k th factor is $\mathbf{C}_{\cdot k}$; has loading $\mathbf{L}_{\cdot k}$ with $p^{-1} \mathbf{L}_{\cdot k}^T \mathbf{L}_{\cdot k} = \lambda_k$. Here λ_k is the average effect size for the k th factor.
- Interpretation: $\mathbf{C}_{\cdot k}$ is the k th (out of K) most important factor.
- Goal: estimate $\mathbf{C}_{\cdot k}$ and $\mathbf{L}_{\cdot k}$.

PCA in noisy data II

$\mathbf{Y}_{p \times n} = \mathbf{L}_{p \times K} \mathbf{C}_{n \times K}^T + \mathbf{E}_{p \times n}$, entries of \mathbf{E} are independent,
 $E_{gi} \sim (0, \sigma_g^2)$

- Claim: there exists a parametrization of \mathbf{L} and \mathbf{C} s.t.

① $E(\mathbf{Y} | \mathbf{C}) = \mathbf{L} \mathbf{C}^T$.

② $n^{-1} \mathbf{C}^T \mathbf{C} = \mathbf{I}_K$ and $p^{-1} \mathbf{L}^T \mathbf{L} = \text{diag}(\lambda_1, \dots, \lambda_K)$,
 $\lambda_1 \geq \dots \geq \lambda_K > 0$.

- To identify $\mathbf{C}_{\cdot 1}, \dots, \mathbf{C}_{\cdot K}$ up to sign, need to assume
 $\lambda_1 > \dots > \lambda_K > 0$. Why?

③ λ_k and $n^{-1/2} \mathbf{C}_{\cdot k}$ are the k th eigenvalue and eigenvector of
 $n^{-1} \mathbf{C} (p^{-1} \mathbf{L}^T \mathbf{L}) \mathbf{C}^T$

- To show this: compute the singular value decomposition of $E(\mathbf{Y} | \mathbf{C}) = \mathbf{L} \mathbf{C}^T$!
- k th factor is $\mathbf{C}_{\cdot k}$; has loading $\mathbf{L}_{\cdot k}$ with $p^{-1} \mathbf{L}_{\cdot k}^T \mathbf{L}_{\cdot k} = \lambda_k$. Here λ_k is the average effect size for the k th factor.
- Interpretation: $\mathbf{C}_{\cdot k}$ is the k th (out of K) most important factor.
- Goal: estimate $\mathbf{C}_{\cdot k}$ and $\mathbf{L}_{\cdot k}$.

PCA in noisy data: estimation

$$\mathbf{Y}_{p \times n} = \mathbf{L}_{p \times K} \mathbf{C}_{n \times K}^T + \mathbf{E}_{p \times n}, \mathbf{E}_{g.} \sim (0, \sigma_g^2 \mathbf{I}_n)$$

- $n^{-1} \mathbf{C}^T \mathbf{C} = \mathbf{I}_K$ and $p^{-1} \mathbf{L}^T \mathbf{L} = \text{diag}(\lambda_1, \dots, \lambda_K)$,
 $\lambda_1 > \dots > \lambda_K > 0$
- λ_k and $n^{-1/2} \mathbf{C}_{.k}$ are the k th eigenvalue and eigenvector of $n^{-1} \mathbf{C} (p^{-1} \mathbf{L}^T \mathbf{L}) \mathbf{C}^T$
- Claim: the first K eigenvectors of $p^{-1} \mathbf{Y}^T \mathbf{Y}$ accurately estimate \mathbf{C} . Idea:
 - 1 $E(p^{-1} \mathbf{Y}^T \mathbf{Y}) = \mathbf{C} (p^{-1} \mathbf{L}^T \mathbf{L}) \mathbf{C}^T + E(p^{-1} \mathbf{E}^T \mathbf{E}) = \mathbf{C} (p^{-1} \mathbf{L}^T \mathbf{L}) \mathbf{C}^T + (p^{-1} \sum_{g=1}^p \sigma_g^2) \mathbf{I}_n = \mathbf{C} (p^{-1} \mathbf{L}^T \mathbf{L}) \mathbf{C}^T + \bar{\sigma}^2 \mathbf{I}_n$
 - 2 This implies eigenvectors of $p^{-1} \mathbf{Y}^T \mathbf{Y} \approx$ eigenvectors of $\mathbf{C} (p^{-1} \mathbf{L}^T \mathbf{L}) \mathbf{C}^T$!
- Note that argument relied on $\mathbf{E}_{g.} \sim (0, \sigma_g^2 \mathbf{I}_n)$. What if samples are related, i.e. $\mathbf{E}_{g.} \sim (0, \mathbf{V}_g)$, $\mathbf{V}_g \neq \sigma_g^2 \mathbf{I}_n$?
- How should we estimate \mathbf{L} ? OLS using estimated design matrix $\hat{\mathbf{C}}$!

PCA in noisy data: estimation

$$\mathbf{Y}_{p \times n} = \mathbf{L}_{p \times K} \mathbf{C}_{n \times K}^T + \mathbf{E}_{p \times n}, \mathbf{E}_{g.} \sim (0, \sigma_g^2 \mathbf{I}_n)$$

- $n^{-1} \mathbf{C}^T \mathbf{C} = \mathbf{I}_K$ and $p^{-1} \mathbf{L}^T \mathbf{L} = \text{diag}(\lambda_1, \dots, \lambda_K)$,
 $\lambda_1 > \dots > \lambda_K > 0$
- λ_k and $n^{-1/2} \mathbf{C}_{.k}$ are the k th eigenvalue and eigenvector of $n^{-1} \mathbf{C} (p^{-1} \mathbf{L}^T \mathbf{L}) \mathbf{C}^T$
- Claim: the first K eigenvectors of $p^{-1} \mathbf{Y}^T \mathbf{Y}$ accurately estimate \mathbf{C} . Idea:

① $E(p^{-1} \mathbf{Y}^T \mathbf{Y}) = \mathbf{C} (p^{-1} \mathbf{L}^T \mathbf{L}) \mathbf{C}^T + E(p^{-1} \mathbf{E}^T \mathbf{E}) =$
 $\mathbf{C} (p^{-1} \mathbf{L}^T \mathbf{L}) \mathbf{C}^T + (p^{-1} \sum_{g=1}^p \sigma_g^2) \mathbf{I}_n = \mathbf{C} (p^{-1} \mathbf{L}^T \mathbf{L}) \mathbf{C}^T + \bar{\sigma}^2 \mathbf{I}_n$

② This implies eigenvectors of $p^{-1} \mathbf{Y}^T \mathbf{Y} \approx$ eigenvectors of $\mathbf{C} (p^{-1} \mathbf{L}^T \mathbf{L}) \mathbf{C}^T$!

- Note that argument relied on $\mathbf{E}_{g.} \sim (0, \sigma_g^2 \mathbf{I}_n)$. What if samples are related, i.e. $\mathbf{E}_{g.} \sim (0, \mathbf{V}_g)$, $\mathbf{V}_g \neq \sigma_g^2 \mathbf{I}_n$?
- How should we estimate \mathbf{L} ? OLS using estimated design matrix $\hat{\mathbf{C}}$!

PCA in noisy data: estimation

$$\mathbf{Y}_{p \times n} = \mathbf{L}_{p \times K} \mathbf{C}_{n \times K}^T + \mathbf{E}_{p \times n}, \mathbf{E}_{g.} \sim (0, \sigma_g^2 \mathbf{I}_n)$$

- $n^{-1} \mathbf{C}^T \mathbf{C} = \mathbf{I}_K$ and $p^{-1} \mathbf{L}^T \mathbf{L} = \text{diag}(\lambda_1, \dots, \lambda_K)$,
 $\lambda_1 > \dots > \lambda_K > 0$
- λ_k and $n^{-1/2} \mathbf{C}_{.k}$ are the k th eigenvalue and eigenvector of $n^{-1} \mathbf{C} (p^{-1} \mathbf{L}^T \mathbf{L}) \mathbf{C}^T$
- Claim: the first K eigenvectors of $p^{-1} \mathbf{Y}^T \mathbf{Y}$ accurately estimate \mathbf{C} . Idea:

① $E(p^{-1} \mathbf{Y}^T \mathbf{Y}) = \mathbf{C} (p^{-1} \mathbf{L}^T \mathbf{L}) \mathbf{C}^T + E(p^{-1} \mathbf{E}^T \mathbf{E}) =$
 $\mathbf{C} (p^{-1} \mathbf{L}^T \mathbf{L}) \mathbf{C}^T + (p^{-1} \sum_{g=1}^p \sigma_g^2) \mathbf{I}_n = \mathbf{C} (p^{-1} \mathbf{L}^T \mathbf{L}) \mathbf{C}^T + \bar{\sigma}^2 \mathbf{I}_n$

② This implies eigenvectors of $p^{-1} \mathbf{Y}^T \mathbf{Y} \approx$ eigenvectors of $\mathbf{C} (p^{-1} \mathbf{L}^T \mathbf{L}) \mathbf{C}^T$!

- Note that argument relied on $\mathbf{E}_{g.} \sim (0, \sigma_g^2 \mathbf{I}_n)$. What if samples are related, i.e. $\mathbf{E}_{g.} \sim (0, \mathbf{V}_g)$, $\mathbf{V}_g \neq \sigma_g^2 \mathbf{I}_n$?
- How should we estimate \mathbf{L} ? OLS using estimated design matrix $\hat{\mathbf{C}}$!

PCA in noisy data: estimation

$$\mathbf{Y}_{p \times n} = \mathbf{L}_{p \times K} \mathbf{C}_{n \times K}^T + \mathbf{E}_{p \times n}, \mathbf{E}_{g.} \sim (0, \sigma_g^2 \mathbf{I}_n)$$

- $n^{-1} \mathbf{C}^T \mathbf{C} = \mathbf{I}_K$ and $p^{-1} \mathbf{L}^T \mathbf{L} = \text{diag}(\lambda_1, \dots, \lambda_K)$,
 $\lambda_1 > \dots > \lambda_K > 0$
- λ_k and $n^{-1/2} \mathbf{C}_{.k}$ are the k th eigenvalue and eigenvector of $n^{-1} \mathbf{C} (p^{-1} \mathbf{L}^T \mathbf{L}) \mathbf{C}^T$
- Claim: the first K eigenvectors of $p^{-1} \mathbf{Y}^T \mathbf{Y}$ accurately estimate \mathbf{C} . Idea:

① $E(p^{-1} \mathbf{Y}^T \mathbf{Y}) = \mathbf{C} (p^{-1} \mathbf{L}^T \mathbf{L}) \mathbf{C}^T + E(p^{-1} \mathbf{E}^T \mathbf{E}) =$
 $\mathbf{C} (p^{-1} \mathbf{L}^T \mathbf{L}) \mathbf{C}^T + (p^{-1} \sum_{g=1}^p \sigma_g^2) \mathbf{I}_n = \mathbf{C} (p^{-1} \mathbf{L}^T \mathbf{L}) \mathbf{C}^T + \bar{\sigma}^2 \mathbf{I}_n$

② This implies eigenvectors of $p^{-1} \mathbf{Y}^T \mathbf{Y} \approx$ eigenvectors of $\mathbf{C} (p^{-1} \mathbf{L}^T \mathbf{L}) \mathbf{C}^T$!

- Note that argument relied on $\mathbf{E}_{g.} \sim (0, \sigma_g^2 \mathbf{I}_n)$. What if samples are related, i.e. $\mathbf{E}_{g.} \sim (0, \mathbf{V}_g)$, $\mathbf{V}_g \neq \sigma_g^2 \mathbf{I}_n$?
- How should we estimate \mathbf{L} ? OLS using estimated design matrix $\hat{\mathbf{C}}$!

PCA in noisy data: estimation

$$\mathbf{Y}_{p \times n} = \mathbf{L}_{p \times K} \mathbf{C}_{n \times K}^T + \mathbf{E}_{p \times n}, \mathbf{E}_{g.} \sim (0, \sigma_g^2 \mathbf{I}_n)$$

- $n^{-1} \mathbf{C}^T \mathbf{C} = \mathbf{I}_K$ and $p^{-1} \mathbf{L}^T \mathbf{L} = \text{diag}(\lambda_1, \dots, \lambda_K)$,
 $\lambda_1 > \dots > \lambda_K > 0$
- λ_k and $n^{-1/2} \mathbf{C}_{.k}$ are the k th eigenvalue and eigenvector of $n^{-1} \mathbf{C} (p^{-1} \mathbf{L}^T \mathbf{L}) \mathbf{C}^T$
- Claim: the first K eigenvectors of $p^{-1} \mathbf{Y}^T \mathbf{Y}$ accurately estimate \mathbf{C} . Idea:

① $E(p^{-1} \mathbf{Y}^T \mathbf{Y}) = \mathbf{C} (p^{-1} \mathbf{L}^T \mathbf{L}) \mathbf{C}^T + E(p^{-1} \mathbf{E}^T \mathbf{E}) =$
 $\mathbf{C} (p^{-1} \mathbf{L}^T \mathbf{L}) \mathbf{C}^T + (p^{-1} \sum_{g=1}^p \sigma_g^2) \mathbf{I}_n = \mathbf{C} (p^{-1} \mathbf{L}^T \mathbf{L}) \mathbf{C}^T + \bar{\sigma}^2 \mathbf{I}_n$

② This implies eigenvectors of $p^{-1} \mathbf{Y}^T \mathbf{Y} \approx$ eigenvectors of $\mathbf{C} (p^{-1} \mathbf{L}^T \mathbf{L}) \mathbf{C}^T$!

- Note that argument relied on $\mathbf{E}_{g.} \sim (0, \sigma_g^2 \mathbf{I}_n)$. What if samples are related, i.e. $\mathbf{E}_{g.} \sim (0, \mathbf{V}_g)$, $\mathbf{V}_g \neq \sigma_g^2 \mathbf{I}_n$?
- How should we estimate \mathbf{L} ? OLS using estimated design matrix $\hat{\mathbf{C}}$!

PCA in noisy data: estimation

$$\mathbf{Y}_{p \times n} = \mathbf{L}_{p \times K} \mathbf{C}_{n \times K}^T + \mathbf{E}_{p \times n}, \mathbf{E}_{g.} \sim (0, \sigma_g^2 \mathbf{I}_n)$$

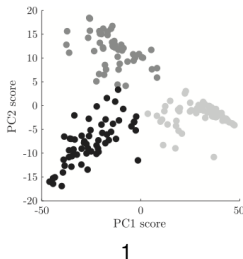
- $n^{-1} \mathbf{C}^T \mathbf{C} = \mathbf{I}_K$ and $p^{-1} \mathbf{L}^T \mathbf{L} = \text{diag}(\lambda_1, \dots, \lambda_K)$,
 $\lambda_1 > \dots > \lambda_K > 0$
- λ_k and $n^{-1/2} \mathbf{C}_{.k}$ are the k th eigenvalue and eigenvector of $n^{-1} \mathbf{C} (p^{-1} \mathbf{L}^T \mathbf{L}) \mathbf{C}^T$
- Claim: the first K eigenvectors of $p^{-1} \mathbf{Y}^T \mathbf{Y}$ accurately estimate \mathbf{C} . Idea:
 - 1 $E(p^{-1} \mathbf{Y}^T \mathbf{Y}) = \mathbf{C} (p^{-1} \mathbf{L}^T \mathbf{L}) \mathbf{C}^T + E(p^{-1} \mathbf{E}^T \mathbf{E}) = \mathbf{C} (p^{-1} \mathbf{L}^T \mathbf{L}) \mathbf{C}^T + (p^{-1} \sum_{g=1}^p \sigma_g^2) \mathbf{I}_n = \mathbf{C} (p^{-1} \mathbf{L}^T \mathbf{L}) \mathbf{C}^T + \bar{\sigma}^2 \mathbf{I}_n$
 - 2 This implies eigenvectors of $p^{-1} \mathbf{Y}^T \mathbf{Y} \approx$ eigenvectors of $\mathbf{C} (p^{-1} \mathbf{L}^T \mathbf{L}) \mathbf{C}^T$!
- Note that argument relied on $\mathbf{E}_{g.} \sim (0, \sigma_g^2 \mathbf{I}_n)$. What if samples are related, i.e. $\mathbf{E}_{g.} \sim (0, \mathbf{V}_g)$, $\mathbf{V}_g \neq \sigma_g^2 \mathbf{I}_n$?
- How should we estimate \mathbf{L} ? OLS using estimated design matrix $\hat{\mathbf{C}}$!

$$\mathbf{Y}_{p \times n} = \mathbf{L}_{p \times K} \mathbf{C}_{n \times K}^T + \mathbf{E}_{p \times n}, \quad \mathbf{E}_{g \cdot} \sim (0, \sigma_g^2 \mathbf{I}_n)$$

- Consider the case $p \gtrsim n$ (we could have $p \gg n$).
- $\text{Corr}(\hat{\mathbf{C}}_{k \cdot}, \mathbf{C}_{k \cdot}) = 1 - O_P \left\{ (\lambda_k np)^{-1/2} + (\lambda_k p)^{-1} \right\}$
 - Blessing of dimensionality: as p gets larger, $\hat{\mathbf{C}}_{k \cdot}$ is more accurate!
 - This is a common theme in factor analysis problems.
- Estimate for ℓ_g is just as accurate as when \mathbf{C} is known (under some assumptions)!
- PCA is incredibly powerful, and can accurately recover factors and loadings.

Example 1: clustering patients based on gene expression

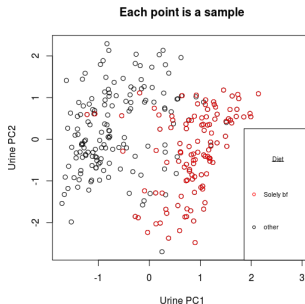
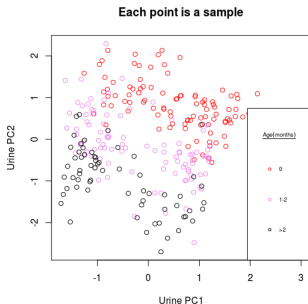
- Expression of $p \approx 15,000$ genes measured on $n = 180$ lung cancer patients.
- Lung cancer has many sub-types (i.e. not just 1 type of lung cancer).
- Goal: can we classify patients based on expression? If we can, this can lead to personalized treatments.



¹Shen et al., 2016

Example 2: identifying confounders in metabolomics

- Metabolomics is the study of small molecule metabolites in tissues/bodily fluids.
 - Represent the end products of all cellular processes.
- What are the major sources of variation?
- Measured concentration of $p = 1,138$ metabolites in the urine of $n = 228$ infants.



2

Example 3: a bit of background

- Recall from high school biology: DNA is made up of the four letter alphabet A, C, T, G.
 - At each locus, inherit 1 copy from mother, and 1 from father
- Single nucleotide polymorphism (SNP): A single base pair (out of $\approx 3 \times 10^9$ base pairs) in the genome that shows variation across populations.
- There are hundreds of millions of SNPs in the human genome.
 - These are random mutations that have been inherited over thousands of generations.
 - At (nearly) every SNP, there is a major (i.e. most frequent) allele and a minor (i.e. less frequent allele).
- Genotype at SNP g in individual i can be written as

$$Y_{gi} = 1 \{i \text{ inherited minor allele at } g \text{ from mother}\} \\ + 1 \{i \text{ inherited minor allele at } g \text{ from father}\}$$

- $Y_{gi} \in \{0, 1, 2\}$

Example 3: a bit of background

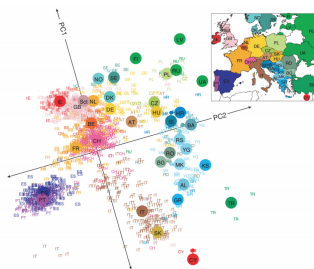
- Recall from high school biology: DNA is made up of the four letter alphabet A, C, T, G.
 - At each locus, inherit 1 copy from mother, and 1 from father
- Single nucleotide polymorphism (SNP): A single base pair (out of $\approx 3 \times 10^9$ base pairs) in the genome that shows variation across populations.
- There are hundreds of millions of SNPs in the human genome.
 - These are random mutations that have been inherited over thousands of generations.
 - At (nearly) every SNP, there is a major (i.e. most frequent) allele and a minor (i.e. less frequent allele).
- Genotype at SNP g in individual i can be written as

$$\begin{aligned} Y_{gi} = & 1 \{i \text{ inherited minor allele at } g \text{ from mother}\} \\ & + 1 \{i \text{ inherited minor allele at } g \text{ from father}\} \end{aligned}$$

- $Y_{gi} \in \{0, 1, 2\}$

Example 3: genes mirror geography!!

- Define $\mathbf{Y} \in \mathbb{R}^{p \times n}$ to be the genotype matrix at $p = 197,146$ SNPs measured in $n = 1,387$ Europeans.
- $Y_{gi} \in \{0, 1, 2\}$ is the genotype.
- Goal: can we cluster individuals based on genotype?
- Intuition: individuals with similar genotypes are more related, and should cluster together.



3

Assessing the angle between subspaces

- Given \mathbf{C} and $\hat{\mathbf{C}}$, how should we assess the angle θ between $\text{im}(\mathbf{C})$ and $\text{im}(\hat{\mathbf{C}})$?
- To define a geometric angle between spaces, they must intersect. Do vector subspaces always intersect?
- When \mathbf{C} and $\hat{\mathbf{C}}$ are vectors, this is easy.
 - $\cos(\theta) = \frac{|\mathbf{C}^T \hat{\mathbf{C}}|}{\|\mathbf{C}\|_2 \|\hat{\mathbf{C}}\|_2}$
- The more general case when $\mathbf{C}, \hat{\mathbf{C}} \in \mathbb{R}^{n \times K}$
 - $\theta = 0 \Leftrightarrow \text{im}(\mathbf{C}) = \text{im}(\hat{\mathbf{C}})$
 - $\theta = \pi/2$ if there exists $\mathbf{v} \in \text{im}(\mathbf{C})$ s.t. \mathbf{v} is orthogonal to $\text{im}(\hat{\mathbf{C}})$.
- $\cos(\theta) = \min_{\hat{\mathbf{v}} \in \text{im}(\hat{\mathbf{C}})} \left\{ \max_{\mathbf{v} \in \text{im}(\mathbf{C})} \left(\frac{\hat{\mathbf{v}}^T \mathbf{v}}{\|\hat{\mathbf{v}}\|_2 \|\mathbf{v}\|_2} \right) \right\}$
 - Also called the first **principal angle** between $\text{im}(\hat{\mathbf{C}})$ and $\text{im}(\mathbf{C})$.

Assessing the angle between subspaces

- Given \mathbf{C} and $\hat{\mathbf{C}}$, how should we assess the angle θ between $\text{im}(\mathbf{C})$ and $\text{im}(\hat{\mathbf{C}})$?
- To define a geometric angle between spaces, they must intersect. Do vector subspaces always intersect?
- When \mathbf{C} and $\hat{\mathbf{C}}$ are vectors, this is easy.
 - $\cos(\theta) = \frac{|\mathbf{C}^T \hat{\mathbf{C}}|}{\|\mathbf{C}\|_2 \|\hat{\mathbf{C}}\|_2}$
- The more general case when $\mathbf{C}, \hat{\mathbf{C}} \in \mathbb{R}^{n \times K}$
 - $\theta = 0 \Leftrightarrow \text{im}(\mathbf{C}) = \text{im}(\hat{\mathbf{C}})$
 - $\theta = \pi/2$ if there exists $\mathbf{v} \in \text{im}(\mathbf{C})$ s.t. \mathbf{v} is orthogonal to $\text{im}(\hat{\mathbf{C}})$.
- $\cos(\theta) = \min_{\hat{\mathbf{v}} \in \text{im}(\hat{\mathbf{C}})} \left\{ \max_{\mathbf{v} \in \text{im}(\mathbf{C})} \left(\frac{\hat{\mathbf{v}}^T \mathbf{v}}{\|\hat{\mathbf{v}}\|_2 \|\mathbf{v}\|_2} \right) \right\}$
 - Also called the first **principal angle** between $\text{im}(\hat{\mathbf{C}})$ and $\text{im}(\mathbf{C})$.

Assessing the angle between subspaces




- Given \mathbf{C} and $\hat{\mathbf{C}}$, how should we assess the angle θ between $\text{im}(\mathbf{C})$ and $\text{im}(\hat{\mathbf{C}})$?
- To define a geometric angle between spaces, they must intersect. Do vector subspaces always intersect?
- When \mathbf{C} and $\hat{\mathbf{C}}$ are vectors, this is easy.
 - $\cos(\theta) = \frac{|\mathbf{C}^T \hat{\mathbf{C}}|}{\|\mathbf{C}\|_2 \|\hat{\mathbf{C}}\|_2}$
- The more general case when $\mathbf{C}, \hat{\mathbf{C}} \in \mathbb{R}^{n \times K}$
 - $\theta = 0 \Leftrightarrow \text{im}(\mathbf{C}) = \text{im}(\hat{\mathbf{C}})$
 - $\theta = \pi/2$ if there exists $\mathbf{v} \in \text{im}(\mathbf{C})$ s.t. \mathbf{v} is orthogonal to $\text{im}(\hat{\mathbf{C}})$.
- $\cos(\theta) = \min_{\hat{\mathbf{v}} \in \text{im}(\hat{\mathbf{C}})} \left\{ \max_{\mathbf{v} \in \text{im}(\mathbf{C})} \left(\frac{\hat{\mathbf{v}}^T \mathbf{v}}{\|\hat{\mathbf{v}}\|_2 \|\mathbf{v}\|_2} \right) \right\}$
 - Also called the first **principal angle** between $\text{im}(\hat{\mathbf{C}})$ and $\text{im}(\mathbf{C})$.

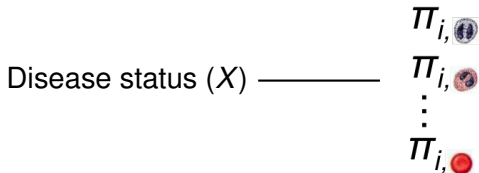
Assessing the angle between subspaces

- Given \mathbf{C} and $\hat{\mathbf{C}}$, how should we assess the angle θ between $\text{im}(\mathbf{C})$ and $\text{im}(\hat{\mathbf{C}})$?
- To define a geometric angle between spaces, they must intersect. Do vector subspaces always intersect?
- When \mathbf{C} and $\hat{\mathbf{C}}$ are vectors, this is easy.
 - $\cos(\theta) = \frac{|\mathbf{C}^T \hat{\mathbf{C}}|}{\|\mathbf{C}\|_2 \|\hat{\mathbf{C}}\|_2}$
- The more general case when $\mathbf{C}, \hat{\mathbf{C}} \in \mathbb{R}^{n \times K}$
 - $\theta = 0 \Leftrightarrow \text{im}(\mathbf{C}) = \text{im}(\hat{\mathbf{C}})$
 - $\theta = \pi/2$ if there exists $\mathbf{v} \in \text{im}(\mathbf{C})$ s.t. \mathbf{v} is orthogonal to $\text{im}(\hat{\mathbf{C}})$.
- $\cos(\theta) = \min_{\hat{\mathbf{v}} \in \text{im}(\hat{\mathbf{C}})} \left\{ \max_{\mathbf{v} \in \text{im}(\mathbf{C})} \left(\frac{\hat{\mathbf{v}}^T \mathbf{v}}{\|\hat{\mathbf{v}}\|_2 \|\mathbf{v}\|_2} \right) \right\}$
 - Also called the first **principal angle** between $\text{im}(\hat{\mathbf{C}})$ and $\text{im}(\mathbf{C})$.

Other applications: adjusting for latent confounders in multivariate regression models

$$Y_i \approx \pi_{i, \text{Cell type 1}} Y_{\text{Cell type 1}} + \pi_{i, \text{Cell type 2}} Y_{\text{Cell type 2}} + \dots + \pi_{i, \text{Cell type } c} Y_{\text{Cell type } c}$$


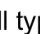

 = Cell type 1,  = Cell type 2, ...,  = Cell type c

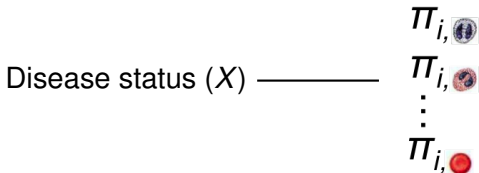


- If we want to understand relationship between disease status and expression, must account for cell type.
- Problem: we only observed expression (Y) and disease status (X). Cell type is unobserved!

Other applications: adjusting for latent confounders in multivariate regression models

$$Y_i \approx \pi_{i, \text{Cell type 1}} Y_{\text{Cell type 1}} + \pi_{i, \text{Cell type 2}} Y_{\text{Cell type 2}} + \dots + \pi_{i, \text{Cell type } c} Y_{\text{Cell type } c}$$

 = Cell type 1,  = Cell type 2, ...,  = Cell type c



- If we want to understand relationship between disease status and expression, must account for cell type.
- Problem: we only observed expression (Y) and disease status (X). Cell type is unobserved!

A model for data with latent confounders

- You now have the tools to model these data!
 - $\mathbf{Y}_g. \in \mathbb{R}^n$: expression at gene g . We measure (i.e. observed) expression.
 - $\mathbf{X} \in \mathbb{R}^n$: covariate of interest, i.e. disease status. This is observed and non-random.
 - $\mathbf{C} \in \mathbb{R}^{n \times K}$: cell type and other latent confounders (e.g. diet, batch).
 - We will ignore other observed nuisance covariates not of interest (e.g. the intercept).
- $\mathbf{Y}_g. = \mathbf{X}\beta_g + \mathbf{C}\ell_g + \mathbf{e}_g, \mathbf{e}_g \sim (\mathbf{0}, \sigma_g^2 I_n)$.
 - Goal: Estimate β_g .
 - Problem: \mathbf{C} may be correlated with \mathbf{X} !
 - Other problem: \mathbf{C} can induce correlations across $\hat{\beta}_1, \dots, \hat{\beta}_p$ (recall Benjamini-Hochberg & other FDR controlling procedures fail when test statistics are correlated).
- We will look at this in depth next time...

A model for data with latent confounders

- You now have the tools to model these data!
 - $\mathbf{Y}_g \in \mathbb{R}^n$: expression at gene g . We measure (i.e. observed) expression.
 - $\mathbf{X} \in \mathbb{R}^n$: covariate of interest, i.e. disease status. This is observed and non-random.
 - $\mathbf{C} \in \mathbb{R}^{n \times K}$: cell type and other latent confounders (e.g. diet, batch).
 - We will ignore other observed nuisance covariates not of interest (e.g. the intercept).
- $\mathbf{Y}_g = \mathbf{X}\beta_g + \mathbf{C}\ell_g + \mathbf{e}_g$, $\mathbf{e}_g \sim (\mathbf{0}, \sigma_g^2 I_n)$.
 - Goal: Estimate β_g .
 - Problem: \mathbf{C} may be correlated with \mathbf{X} !
 - Other problem: \mathbf{C} can induce correlations across $\hat{\beta}_1, \dots, \hat{\beta}_p$ (recall Benjamini-Hochberg & other FDR controlling procedures fail when test statistics are correlated).
- We will look at this in depth next time...