

Applied Statistical Methods II

Some additional topics II

How do we choose the number of latent factors

- $\mathbf{Y}_{p \times n} = \mathbf{L}_{p \times K} \mathbf{C}_{n \times K}^T + \mathbf{E}_{p \times n}$
- The rows $g = 1, \dots, p$ of \mathbf{Y} are genes (or survey questions, companies, etc.).
 - We assume each row behaves as a standard linear model, $\mathbf{E}_{g \cdot} \sim (\mathbf{0}, \sigma^2 \mathbf{I}_n)$
 - Entries of \mathbf{E} are independent.
- The columns $i = 1, \dots, n$ of \mathbf{Y} are samples (or individuals, time, etc.)
- Ultimate goal is to understand $\mathbf{L} \in \mathbb{R}^{p \times K}$ and $\mathbf{C} \in \mathbb{R}^{n \times K}$.
- Problem: we don't know K ! How should we estimate it?
- Choosing K is likely the most important part of exploratory factor analysis (Brown, 2014).

How do we choose the number of latent factors

- $\mathbf{Y}_{p \times n} = \mathbf{L}_{p \times K} \mathbf{C}_{n \times K}^T + \mathbf{E}_{p \times n}$
- The rows $g = 1, \dots, p$ of \mathbf{Y} are genes (or survey questions, companies, etc.).
 - We assume each row behaves as a standard linear model, $\mathbf{E}_{g \cdot} \sim (\mathbf{0}, \sigma^2 \mathbf{I}_n)$
 - Entries of \mathbf{E} are independent.
- The columns $i = 1, \dots, n$ of \mathbf{Y} are samples (or individuals, time, etc.)
- Ultimate goal is to understand $\mathbf{L} \in \mathbb{R}^{p \times K}$ and $\mathbf{C} \in \mathbb{R}^{n \times K}$.
- Problem: we don't know K ! How should we estimate it?
- Choosing K is likely the most important part of exploratory factor analysis (Brown, 2014).

Using a scree plot

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

- Plot the eigenvalues of $p^{-1} \mathbf{Y}^T \mathbf{Y}$.
- Recall

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

- Idea: if the eigenvalues of the signal $\mathbf{C} (p^{-1} \mathbf{L}^T \mathbf{L}) \mathbf{C}^T$ are large, they should be much bigger than those from the noise $p^{-1} \mathbf{E}^T \mathbf{E}$.
- Look for an **elbow** in the scree plot: a point below which the eigenvalues are “small” and can be ignored.
- Problems: qualitative, usual does not work in real data.

Using a scree plot

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

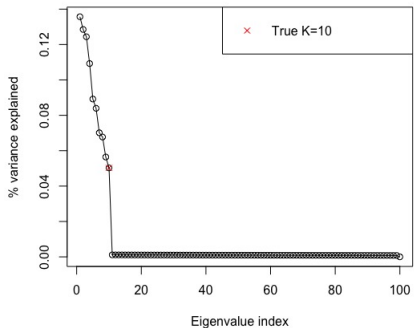
- Plot the eigenvalues of $p^{-1} \mathbf{Y}^T \mathbf{Y}$.
- Recall

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

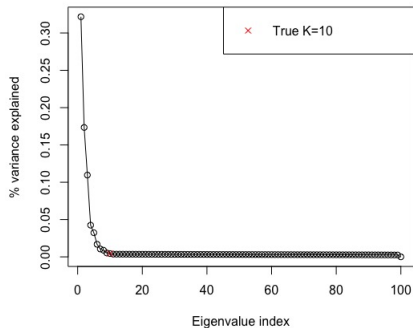
- Idea: if the eigenvalues of the signal $\mathbf{C} (p^{-1} \mathbf{L}^T \mathbf{L}) \mathbf{C}^T$ are large, they should be much bigger than those from the noise $p^{-1} \mathbf{E}^T \mathbf{E}$.
- Look for an **elbow** in the scree plot: a point below which the eigenvalues are “small” and can be ignored.
- Problems: qualitative, usual does not work in real data.

Using a scree plot: simulated data

Simulated data, large eigenvalues

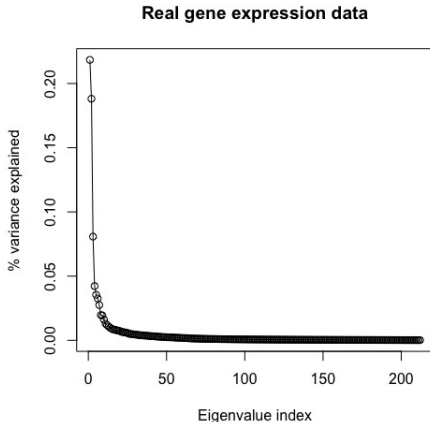


Simulated data, small eigenvalues



Using a scree plot: real gene expression data

- Real gene expression data from Knowles et al., 2018.
- Measured the expression of $p = 12,317$ genes in $n = 217$ samples.
- What value of K would you use?



More quantitative approaches to choose K

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

- Bai & Ng (2002): A BIC-like criteria.
 - $f(k) = R(k) + kP(p, n)$. Choose k to **minimize** $f(k)$.
 - $R(k) = (np)^{-1} \|\mathbf{Y} - \hat{\mathbf{L}}^{(k)} \left\{ \hat{\mathbf{C}}^{(k)} \right\}^T\|_F^2$ are the sum of squared residuals for PCA's estimates when it is assumed $K = k$.
 - $P(p, n)$ is a penalty function that satisfies $P(p, n) \rightarrow 0$ and $\min(n, p)P(p, n) \rightarrow \infty$ as $n, p \rightarrow \infty$.
 - Example: $P(p, n) = 1 / \log \{\min(n, p)\}$
 - This is a groundbreaking paper.
 - **Serious problem:** This only works in theory and in practice if the eigenvalues of $\mathbf{C} (p^{-1} \mathbf{L}^T \mathbf{L}) \mathbf{C}^T$ are **very large** (i.e. $\asymp n$).
 - In these cases, we might as well use a scree plot!
- There are lots of similar methods (Ahn & Horenstein, 2013; Lu & Su, 2016; Li et al., 2018).
- All suffer from the same problems.

More quantitative approaches to choose K

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

- Bai & Ng (2002): A BIC-like criteria.
 - $f(k) = R(k) + kP(p, n)$. Choose k to **minimize** $f(k)$.
 - $R(k) = (np)^{-1} \|\mathbf{Y} - \hat{\mathbf{L}}^{(k)} \{\hat{\mathbf{C}}^{(k)}\}^T\|_F^2$ are the sum of squared residuals for PCA's estimates when it is assumed $K = k$.
 - $P(p, n)$ is a penalty function that satisfies $P(p, n) \rightarrow 0$ and $\min(n, p)P(p, n) \rightarrow \infty$ as $n, p \rightarrow \infty$.
 - Example: $P(p, n) = 1 / \log \{\min(n, p)\}$
 - This is a groundbreaking paper.
 - **Serious problem:** This only works in theory and in practice if the eigenvalues of $\mathbf{C} (p^{-1} \mathbf{L}^T \mathbf{L}) \mathbf{C}^T$ are **very large** (i.e. $\asymp n$).
 - In these cases, we might as well use a scree plot!
- There are lots of similar methods (Ahn & Horenstein, 2013; Lu & Su, 2016; Li et al., 2018).
- All suffer from the same problems.

More quantitative approaches to choose K

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

- Bai & Ng (2002): A BIC-like criteria.
 - $f(k) = R(k) + kP(p, n)$. Choose k to **minimize** $f(k)$.
 - $R(k) = (np)^{-1} \|\mathbf{Y} - \hat{\mathbf{L}}^{(k)} \left\{ \hat{\mathbf{C}}^{(k)} \right\}^T\|_F^2$ are the sum of squared residuals for PCA's estimates when it is assumed $K = k$.
 - $P(p, n)$ is a penalty function that satisfies $P(p, n) \rightarrow 0$ and $\min(n, p)P(p, n) \rightarrow \infty$ as $n, p \rightarrow \infty$.
 - Example: $P(p, n) = 1 / \log \{\min(n, p)\}$
 - This is a groundbreaking paper.
 - **Serious problem:** This only works in theory and in practice if the eigenvalues of $\mathbf{C} (p^{-1} \mathbf{L}^T \mathbf{L}) \mathbf{C}^T$ are **very large** (i.e. $\asymp n$).
 - In these cases, we might as well use a scree plot!
- There are lots of similar methods (Ahn & Horenstein, 2013; Lu & Su, 2016; Li et al., 2018).
- All suffer from the same problems.

More quantitative approaches to choose K

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

- Bai & Ng (2002): A BIC-like criteria.
 - $f(k) = R(k) + kP(p, n)$. Choose k to **minimize** $f(k)$.
 - $R(k) = (np)^{-1} \|\mathbf{Y} - \hat{\mathbf{L}}^{(k)} \left\{ \hat{\mathbf{C}}^{(k)} \right\}^T\|_F^2$ are the sum of squared residuals for PCA's estimates when it is assumed $K = k$.
 - $P(p, n)$ is a penalty function that satisfies $P(p, n) \rightarrow 0$ and $\min(n, p)P(p, n) \rightarrow \infty$ as $n, p \rightarrow \infty$.
 - Example: $P(p, n) = 1 / \log \{\min(n, p)\}$
 - This is a groundbreaking paper.
 - **Serious problem:** This only works in theory and in practice if the eigenvalues of $\mathbf{C} (p^{-1} \mathbf{L}^T \mathbf{L}) \mathbf{C}^T$ are **very large** (i.e. $\asymp n$).
 - In these cases, we might as well use a scree plot!
- There are lots of similar methods (Ahn & Horenstein, 2013; Lu & Su, 2016; Li et al., 2018).
- All suffer from the same problems.

Parallel analysis (PA) to choose K

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

- A more data-driven approach, based on permutation.
- $\mathbf{Y} = \mathbf{Signal} + \mathbf{Noise}$.
- Idea: we only include a factor k if it's eigenvalue is suitably greater than that of the noise.
- Goal: need to understand singular values of the noise.
- **Option 1:** use random matrix theory (RMT) to understand the singular values of \mathbf{E} .
 - There are some really beautiful results here. See Chapter 5 of Eldar & Kutyniok (2012) for some of them.
 - Problem: results are sensitive to distributional assumptions on \mathbf{E} .
- **Option 2:** Use the observed data to estimate singular values of \mathbf{E} .
 - Called parallel analysis (PA).
 - A very old idea (dates back to at least 1992). The first theoretical results were recently published (Dobriban, 2020).

Parallel analysis (PA) to choose K

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

- A more data-driven approach, based on permutation.
- $\mathbf{Y} = \mathbf{Signal} + \mathbf{Noise}$.
- Idea: we only include a factor k if it's eigenvalue is suitably greater than that of the noise.
- Goal: need to understand singular values of the noise.
- **Option 1:** use random matrix theory (RMT) to understand the singular values of \mathbf{E} .
 - There are some really beautiful results here. See Chapter 5 of Eldar & Kutyniok (2012) for some of them.
 - Problem: results are sensitive to distributional assumptions on \mathbf{E} .
- **Option 2:** Use the observed data to estimate singular values of \mathbf{E} .
 - Called parallel analysis (PA).
 - A very old idea (dates back to at least 1992). The first theoretical results were recently published (Dobriban, 2020).

Parallel analysis (PA) to choose K

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

- A more data-driven approach, based on permutation.
- $\mathbf{Y} = \mathbf{Signal} + \mathbf{Noise}$.
- Idea: we only include a factor k if it's eigenvalue is suitably greater than that of the noise.
- Goal: need to understand singular values of the noise.
- **Option 1:** use random matrix theory (RMT) to understand the singular values of \mathbf{E} .
 - There are some really beautiful results here. See Chapter 5 of Eldar & Kutyniok (2012) for some of them.
 - Problem: results are sensitive to distributional assumptions on \mathbf{E} .
- **Option 2:** Use the observed data to estimate singular values of \mathbf{E} .
 - Called parallel analysis (PA).
 - A very old idea (dates back to at least 1992). The first theoretical results were recently published (Dobriban, 2020).

PA: the algorithm

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

- Mean center the rows of \mathbf{Y} , i.e. $\mathbf{Y} \leftarrow \mathbf{Y} (\mathbf{I}_n - n^{-1} \mathbf{1}_n \mathbf{1}_n^T)$
- Let $\delta_1, \delta_2, \dots, \delta_{\min(n,p)}$ be the singular values of \mathbf{Y} .
- Idea: choose \hat{K} s.t. $\delta_{\hat{K}}$ is suitably larger than singular values of the noise \mathbf{E} .
- To approximate the singular values of the noise \mathbf{E} :
 - 1 “Break” the signal $\mathbf{L}\mathbf{C}^T$ by independently permuting entries in the rows of \mathbf{Y} : $\tilde{\mathbf{E}}_{g.}^{(b)} = \Pi_g \mathbf{Y}_{g.}$, Π_g a random permutation matrix.
 - 2 Let $\delta_1^{(b)}, \delta_2^{(b)}, \dots, \delta_{\min(n,p)}^{(b)}$ be the s.v. of $\tilde{\mathbf{E}}^{(b)}$.
 - 3 Repeat for $b = 1, \dots, B$. This gives you an approximate dist’n of the σ_k (\mathbf{E}) (the k th singular value of \mathbf{E}).
- Keep factor k if $\delta_k > 95\%$ of the of the $\delta_k^{(b)}$ ’s.
- **Problem:** $\tilde{\mathbf{E}}^{(b)}$ is often very different from \mathbf{E} .
- Subject to **eigenvalue shadowing**: $\tilde{\mathbf{E}}^{(b)}$ is contaminated by factors with large eigenvalues.

PA: the algorithm

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

- Mean center the rows of \mathbf{Y} , i.e. $\mathbf{Y} \leftarrow \mathbf{Y} \left(\mathbf{I}_n - n^{-1} \mathbf{1}_n \mathbf{1}_n^T \right)$
- Let $\delta_1, \delta_2, \dots, \delta_{\min(n,p)}$ be the singular values of \mathbf{Y} .
- Idea: choose \hat{K} s.t. $\delta_{\hat{K}}$ is suitably larger than singular values of the noise \mathbf{E} .
- To approximate the singular values of the noise \mathbf{E} :
 - 1 “Break” the signal $\mathbf{L}\mathbf{C}^T$ by independently permuting entries in the rows of \mathbf{Y} : $\tilde{\mathbf{E}}_g^{(b)} = \mathbf{\Pi}_g \mathbf{Y}_g$, $\mathbf{\Pi}_g$ a random permutation matrix.
 - 2 Let $\delta_1^{(b)}, \delta_2^{(b)}, \dots, \delta_{\min(n,p)}^{(b)}$ be the s.v. of $\tilde{\mathbf{E}}^{(b)}$.
 - 3 Repeat for $b = 1, \dots, B$. This gives you an approximate dist’n of the σ_k (\mathbf{E}) (the k th singular value of \mathbf{E}).
- Keep factor k if $\delta_k > 95\%$ of the of the $\delta_k^{(b)}$ ’s.
- **Problem:** $\tilde{\mathbf{E}}^{(b)}$ is often very different from \mathbf{E} .
- Subject to **eigenvalue shadowing**: $\tilde{\mathbf{E}}^{(b)}$ is contaminated by factors with large eigenvalues.

PA: the algorithm

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

- Mean center the rows of \mathbf{Y} , i.e. $\mathbf{Y} \leftarrow \mathbf{Y} \left(\mathbf{I}_n - n^{-1} \mathbf{1}_n \mathbf{1}_n^T \right)$
- Let $\delta_1, \delta_2, \dots, \delta_{\min(n,p)}$ be the singular values of \mathbf{Y} .
- Idea: choose \hat{K} s.t. $\delta_{\hat{K}}$ is suitably larger than singular values of the noise \mathbf{E} .
- To approximate the singular values of the noise \mathbf{E} :
 - 1 “Break” the signal $\mathbf{L}\mathbf{C}^T$ by independently permuting entries in the rows of \mathbf{Y} : $\tilde{\mathbf{E}}_g^{(b)} = \mathbf{\Pi}_g \mathbf{Y}_g$, $\mathbf{\Pi}_g$ a random permutation matrix.
 - 2 Let $\delta_1^{(b)}, \delta_2^{(b)}, \dots, \delta_{\min(n,p)}^{(b)}$ be the s.v. of $\tilde{\mathbf{E}}^{(b)}$.
 - 3 Repeat for $b = 1, \dots, B$. This gives you an approximate dist’n of the σ_k (\mathbf{E}) (the k th singular value of \mathbf{E}).
- Keep factor k if $\delta_k > 95\%$ of the of the $\delta_k^{(b)}$ ’s.
- **Problem:** $\tilde{\mathbf{E}}^{(b)}$ is often very different from \mathbf{E} .
- Subject to **eigenvalue shadowing**: $\tilde{\mathbf{E}}^{(b)}$ is contaminated by factors with large eigenvalues.

PA: the algorithm

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

- Mean center the rows of \mathbf{Y} , i.e. $\mathbf{Y} \leftarrow \mathbf{Y} \left(\mathbf{I}_n - n^{-1} \mathbf{1}_n \mathbf{1}_n^T \right)$
- Let $\delta_1, \delta_2, \dots, \delta_{\min(n,p)}$ be the singular values of \mathbf{Y} .
- Idea: choose \hat{K} s.t. $\delta_{\hat{K}}$ is suitably larger than singular values of the noise \mathbf{E} .
- To approximate the singular values of the noise \mathbf{E} :
 - 1 “Break” the signal $\mathbf{L}\mathbf{C}^T$ by independently permuting entries in the rows of \mathbf{Y} : $\tilde{\mathbf{E}}_g^{(b)} = \mathbf{\Pi}_g \mathbf{Y}_g$, $\mathbf{\Pi}_g$ a random permutation matrix.
 - 2 Let $\delta_1^{(b)}, \delta_2^{(b)}, \dots, \delta_{\min(n,p)}^{(b)}$ be the s.v. of $\tilde{\mathbf{E}}^{(b)}$.
 - 3 Repeat for $b = 1, \dots, B$. This gives you an approximate dist’n of the σ_k (\mathbf{E}) (the k th singular value of \mathbf{E}).
- Keep factor k if $\delta_k > 95\%$ of the of the $\delta_k^{(b)}$ ’s.
- **Problem:** $\tilde{\mathbf{E}}^{(b)}$ is often very different from \mathbf{E} .
- Subject to **eigenvalue shadowing**: $\tilde{\mathbf{E}}^{(b)}$ is contaminated by factors with large eigenvalues.

Other algorithms

$\mathbf{Y}_{p \times n} = \mathbf{L}_{p \times K} \mathbf{C}_{n \times K}^T + \mathbf{E}_{p \times n}$, $\mathbf{E}_g \sim (\mathbf{0}, \sigma^2 \mathbf{I}_n)$. Assume rows of \mathbf{E} are independent.

- Bi-cross validation (Wang & Owen, 2016; McKennan & Nicolae, 2019).
- Idea: partition columns of \mathbf{Y} into training (f) and test ($-f$) sets

$$\mathbf{Y} = \begin{bmatrix} \mathbf{Y}_{(-f)} \\ \mathbf{Y}_f \end{bmatrix} = \begin{bmatrix} \mathbf{L}_{(-f)} \mathbf{C}^T \\ \mathbf{L}_f \mathbf{C}^T \end{bmatrix} + \begin{bmatrix} \mathbf{E}_{(-f)} \\ \mathbf{E}_f \end{bmatrix}$$

- Key observation: $\mathbf{E}_{(-f)}$ is independent of \mathbf{E}_f !
- Can estimate \mathbf{C} from \mathbf{Y}_f .
- Test estimate using $\mathbf{Y}_{(-f)}$.
- In my (biased) opinion, this is the most reliable way to estimate K .
- Has really beautiful theory (McKennan & Nicolae, 2019).