# Efficiency of the CPC estimator in modelling the covariance structures of several populations

Theo Pepler
Unit for Biometry
Stellenbosch University

**Principal component analysis** (PCA):

$$\Sigma = \beta \Lambda \beta'$$

**Common principal components** (CPC):

$$\Sigma_1 = \beta \Lambda_1 \beta'$$

$$\Sigma_2 = \beta \Lambda_2 \beta'$$

Estimate $\beta$ with Flury-Gautschi (or other) algorithm.

**Question:** If the CPC hypothesis is tenable, can the information about the common eigenvectors be used to find improved estimates of the covariance matrices of the populations?

**CPC estimator of $\Sigma_i$ (Flury, 1988)**

- $S_i$ : unbiased sample covariance matrix estimator
- $B$ : estimator of modal matrix, $\beta$

$$L_i = B'S_iB \tag{1}$$

$$L_i^0 = \text{diag}(L_i) \tag{2}$$

$$S_{i(CPC)} = BL_i^0B' \tag{3}$$

### Regularised estimator of $\Sigma_i$ (Friedman, 1989)

- $S_{\text{pool}}$ : pooled covariance matrix estimator
- $\alpha_i$ : shrinkage intensity parameter

$$S_i^\star = \alpha_i S_i + (1 - \alpha_i) S_{\text{pool}} \tag{4}$$

Substitute $S_{i(CPC)}$ for $S_{\text{pool}}$ in Equation (4) to find regularised CPC estimator:

$$S_{i(CPC)}^\star = \alpha_i S_i + (1 - \alpha_i) S_{i(CPC)} \tag{5}$$

Substitute $S_i = BL_iB$ and $S_{i(CPC)} = BL_i^0B$ in Equation (5):

$$S_{i(\text{CPC})}^\star = \alpha_i BL_iB' + (1 - \alpha_i)BL_i^0B'$$

$$= B[\alpha_i L_i + (1 - \alpha_i)L_i^0]B'$$

$$= B[\alpha_i L_i + L_i^0 - \alpha_i L_i^0]B' \quad (6)$$

$$= B[\alpha_i(L_i - L_i^0) + L_i^0]B'$$

**Estimation of $\alpha_i$:**

- inverse of $\phi$ measure (Flury, 1988):

$$\hat{\alpha}_i = 1 - \phi(\boldsymbol{L}_i)^{-1}$$

$$= 1 - \frac{\det(\boldsymbol{L}_i)}{\det\left[\text{diag}(\boldsymbol{L}_i)\right]} \qquad (7)$$

$$= 1 - \frac{\det(\boldsymbol{L}_i)}{\det(\boldsymbol{L}_i^0)}$$

**Estimation of $\alpha_i$:**

- optimisation on validation data
  Group $i$ sample

  $r = 1, \ldots, 100$ replications

| | |
|---|---|
| **70%** | $\rightarrow \boldsymbol{S}_{i(TRAIN)}^{(r)}, \boldsymbol{S}_{i(CPC)}^{(r)}$ |
| **30%** | $\rightarrow \boldsymbol{S}_{i(VALID)}^{(r)}$ <br> Find $\alpha_i^{(r)}$ which minimises <br> $\left\| \left[ \alpha_i^{(r)} \boldsymbol{S}_{i(TRAIN)}^{(r)} + (1 - \alpha_i^{(r)}) \boldsymbol{S}_{i(CPC)}^{(r)} \right] - \boldsymbol{S}_{i(VALID)}^{(r)} \right\|_{F^\star}$ |

$$\hat{\alpha}_i = \frac{\sum_r \alpha_i^{(r)}}{r} \tag{8}$$

**Estimation of $\alpha_i$:**

- method adapted from Schäfer & Strimmer (2005):

$$\hat{\alpha}_i = 1 - \min\left(\frac{\sum_{j\neq h}\hat{\mathsf{Var}}(l_{ijh})}{\sum_{j\neq h} l_{ijh}^2}, 1\right) \qquad (9)$$

$l_{ijh}$ : element in the $j^{th}$ row and $h^{th}$ column of $\boldsymbol{L}_i$

Estimate $\hat{\mathsf{Var}}(l_{ijh})$ with bootstrap.

**Monte Carlo simulation comparing estimators:**

**1** Unbiased sample covariance matrix

**2** CPC estimator

**3** Regularised CPC estimator
- inverse $\phi$
- optimisation with validation data
- Schäfer & Strimmer method with bootstrap
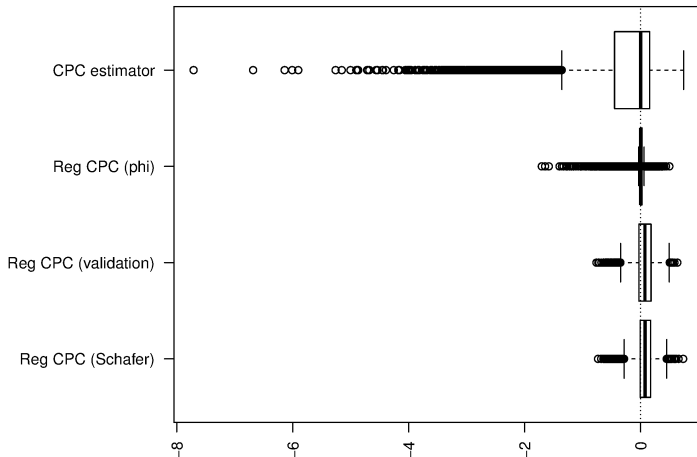
**Simulation parameters**

- $k = 2$ groups

- $p = 5, 10, 20$ variables

- Common eigenvector rank orders: same, similar, opposite

- Group 1 sample size: $n_1 = 200, 500, 1000$

- Group 2 sample size: $n_2 = 30, 50, 100, 200$

- Multivariate distributions: normal, chi-square (2 df), $t$ (1 df)
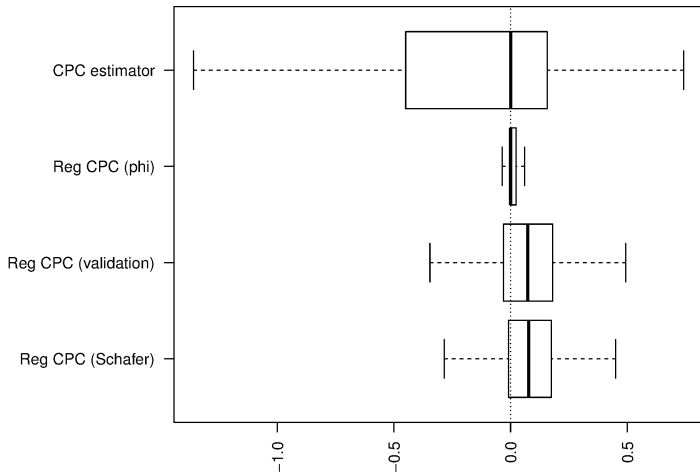
**Modified Frobenius norm**

Error of estimation when comparing $\hat{\mathbf{\Sigma}}_i$ to $\mathbf{\Sigma}_i$:

$$||\hat{\mathbf{\Sigma}}_i - \mathbf{\Sigma}_i||_{F^\star} = \sqrt{\sum_{j=1}^{p}\sum_{h \geq j}^{p}(\hat{\sigma}_{ijh} - \sigma_{ijh})^2} \qquad (10)$$
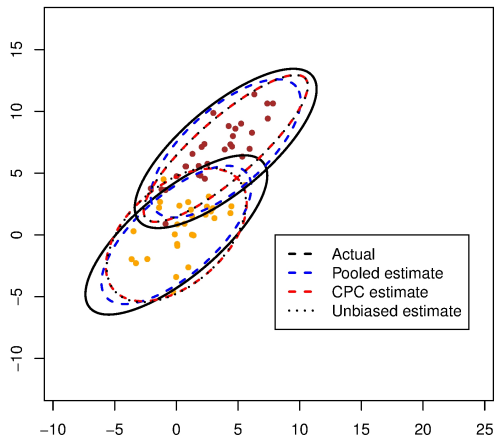
Improvement in modified Frobenius norm
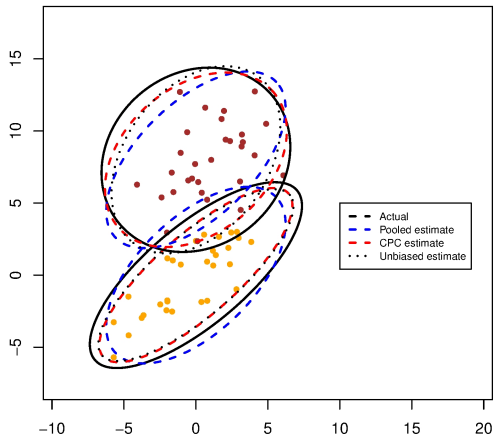
Improvement in modified Frobenius norm

**Conclusions**

- Regularised CPC estimator ($\alpha_i$ found with validation method) performs best

- Median improvement in modified Frobenius norm of up to 27% (Multivariate $t$ data, $n_1 = 1000$, $n_2 = 30$)

- Using CPC model in covariance matrix estimation provides greatest benefit for groups with *small $n_i : p$ ratios*

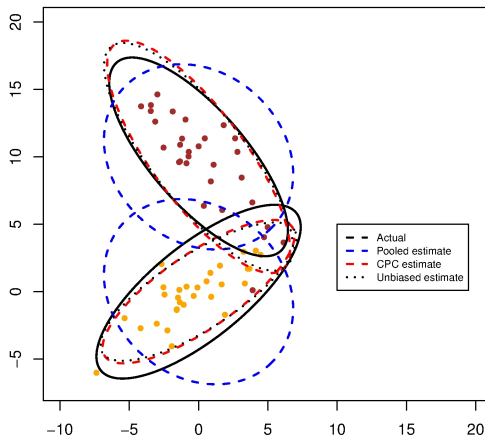- Regularised CPC estimators offered greater improvement (over unbiased estimator) with *non-normal data*
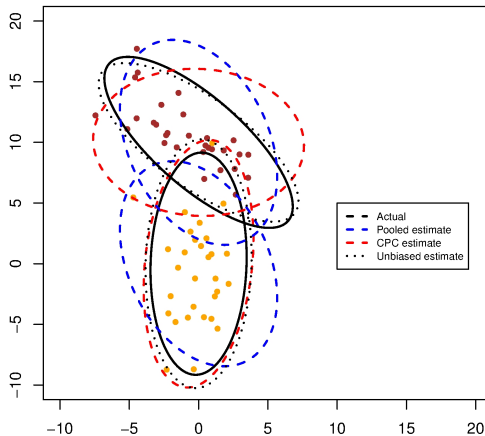
**Equal covariance matrices**

Same common eigenvector rank order

Opposite common eigenvector rank order

**Unrelated covariance matrices**

Legend:
- Actual
- Pooled estimate
- CPC estimate
- Unbiased estimate

**References**

Flury, B. (1988). *Common principal components and related multivariate models*.

Friedman, J.H. (1989). Regularized discriminant analysis. *Journal of the American Statistical Association*, **405**, 165–175.

Schäfer, J., Strimmer, K. (2005). A shrinkage approach to large-scale covariance matrix estimation and implications for functional genomics. *Statistical applications in genetics and molecular biology*, **4:1**, 1175-1189.