

Dimension Reduction with Heavy Tails

Gabriel Kuhn

Munich University of Technology

<http://www.ma.tum.de/stat>

4th Conference on Extreme Value Analysis

Gothenburg, August 15.–19., 2005

Reference: Klüppelberg, C. and Kuhn, G. (2005) Dimension Reduction with Heavy Tails. *In preparation.*

Factor Model

- Observable d -dimensional random vector \mathbf{X}
- **Model:** $\mathbf{X} = \boldsymbol{\mu} + \mathbf{L}\mathbf{f} + \mathbf{V}\mathbf{e}$
 - $\mathbf{L}\mathbf{f}$: k -dimensional non-observable *common factors* \mathbf{f} , loading matrix \mathbf{L}
 - $\mathbf{V}\mathbf{e}$: *specific factors* \mathbf{e} , diagonal matrix \mathbf{V}
 - (\mathbf{f}, \mathbf{e}) are uncorrelated (independent)
- **Idea:** Distribution of \mathbf{X} described by linear combination of k factors with componentwise extra source of randomness.
- **Classical model:** $(\mathbf{f}, \mathbf{e}) \sim \mathcal{N}_{k+d}(\mathbf{0}, \mathbf{I})$ and $\text{Cov}(\mathbf{X}) =: \boldsymbol{\Sigma} = \mathbf{L}\mathbf{L}^T + \mathbf{V}^2$
- **Disadvantages:**
 - Data may not be normal
 - No heavy-tailed model
 - Dependence in extremes cannot be modeled
 - Margins of the same type
- **Task:** Overcome disadvantages above

Distribution Models

- *Elliptical Distribution:* $\mathbf{X} \stackrel{d}{=} \boldsymbol{\mu} + G\mathbf{A}\mathbf{U}^{(k)}$
 $G > 0$ independent of $\mathbf{U}^{(k)} \sim \text{unif}\{\mathbf{s} : \|\mathbf{s}\| = 1\}$, $\mathbf{A} \in \mathbb{R}^{d \times k}$, $\boldsymbol{\Sigma} := \mathbf{A}\mathbf{A}^T$
 $E(\mathbf{X}) = \boldsymbol{\mu}$, $\text{Cov}(\mathbf{X}) = EG^2\boldsymbol{\Sigma}/k$, $\text{Corr}(\mathbf{X}) = \text{diag}(\boldsymbol{\Sigma})^{-1/2}\boldsymbol{\Sigma}\text{diag}(\boldsymbol{\Sigma})^{-1/2} =: \mathbf{R}$
- **Example:**
 - *normal:* $\mathbf{X} \stackrel{d}{=} \boldsymbol{\mu} + \sqrt{\chi_k^2}\mathbf{A}\mathbf{U}^{(k)} \sim \mathcal{N}_d(\boldsymbol{\mu}, \boldsymbol{\Sigma})$
 - *multivariate t_ν :* $\mathbf{X} \stackrel{d}{=} \boldsymbol{\mu} + \sqrt{\nu\chi_k^2/\chi_\nu^2}\mathbf{A}\mathbf{U}^{(k)} \stackrel{d}{=} \boldsymbol{\mu} + \sqrt{\nu/\chi_\nu^2}\mathcal{N}_d(\mathbf{0}, \boldsymbol{\Sigma})$

Extended Factor Model

Drop normal assumption and consider (with $\mathbf{L}\mathbf{L}^T + \mathbf{V}^2 = \Sigma$):

$$\mathbf{X} \stackrel{d}{=} \boldsymbol{\mu} + \mathbf{L}\mathbf{f} + \mathbf{V}\mathbf{e} \text{ is elliptical: } \mathbf{L}\mathbf{f} + \mathbf{V}\mathbf{e} = (\mathbf{L}, \mathbf{V}) \begin{pmatrix} \mathbf{f} \\ \mathbf{e} \end{pmatrix} \stackrel{d}{=} G(\mathbf{L}, \mathbf{V})U^{(k+d)},$$

e.g. choose multivariate t_ν -distribution

- Remark:**
- \mathbf{f} and \mathbf{e} are uncorrelated but not independent
 - alternatively: same dependence structure, but arbitrary margins (copula approach)
 - $P(G > x) \sim Cx^{-\nu}$ needed for modelling tail dependence

Factorization

- Standard approach uses (normal) ml algorithm for decomposition $\Sigma = \mathbf{L}\mathbf{L}^T + \mathbf{V}^2$

Definition: $f_{\text{ml}}^{(k)}(\mathbf{A}) = (\mathbf{L}, \mathbf{V}), \quad (f_{\text{ml}}^{(k)}(\mathbf{A}))^2 := (\mathbf{L}, \mathbf{V})(\mathbf{L}, \mathbf{V})^T = \mathbf{L}\mathbf{L}^T + \mathbf{V}^2$

- **Lemma:** $\Sigma_n \xrightarrow{P} \Sigma = \mathbf{L}\mathbf{L}^T + \mathbf{V}^2$ and neighborhood of Σ decomposable by $f_{\text{ml}}^{(k)}$
 $\Rightarrow \left(f_{\text{ml}}^{(k)}(\Sigma_n) \right)^2 \xrightarrow{P} \Sigma$

- **Interpretation:** Given some *consistent* and *composable* covariance (correlation) estimator, the algorithm computes a consistent decomposition
(independent of distribution model)
- **Remark:** In application algorithm almost always produces a decomposition

Dependence Concepts

- *Kendall's τ* : Let $(X, Y) \stackrel{iid}{\sim} (\tilde{X}, \tilde{Y})$

$$\tau := P\left((X - \tilde{X})(Y - \tilde{Y}) > 0\right) - P\left((X - \tilde{X})(Y - \tilde{Y}) < 0\right)$$

Elliptical distribution $\Rightarrow \mathbf{R} = \sin(\pi \mathbf{T} / 2), \quad \mathbf{T} = (\tau_{ij})_{1 \leq i, j \leq d}$

- *Tail Dependence*: (X, Y) with margins F_X, F_Y

$$\lambda := \lim_{u \searrow 0} P(Y < F_Y^{\leftarrow}(u) \mid X < F_X^{\leftarrow}(u))$$

Elliptical distribution and $P(G > x) \sim Cx^{-\nu}, \nu \in (0, \infty)$

$$\Rightarrow \mathbf{R} = 1 - 2 \frac{(F_{t, \nu+1}^{\leftarrow}(1 - \Lambda/2))^2}{\nu + 1 + (F_{t, \nu+1}^{\leftarrow}(1 - \Lambda/2))^2}$$

$$\mathbf{\Lambda} = (\lambda_{ij})_{1 \leq i, j \leq d} \text{ and } F_{t, \nu} : \text{df of 1-dim } t_\nu$$

- *Remark*: Both dependence concepts independent of margins

Example

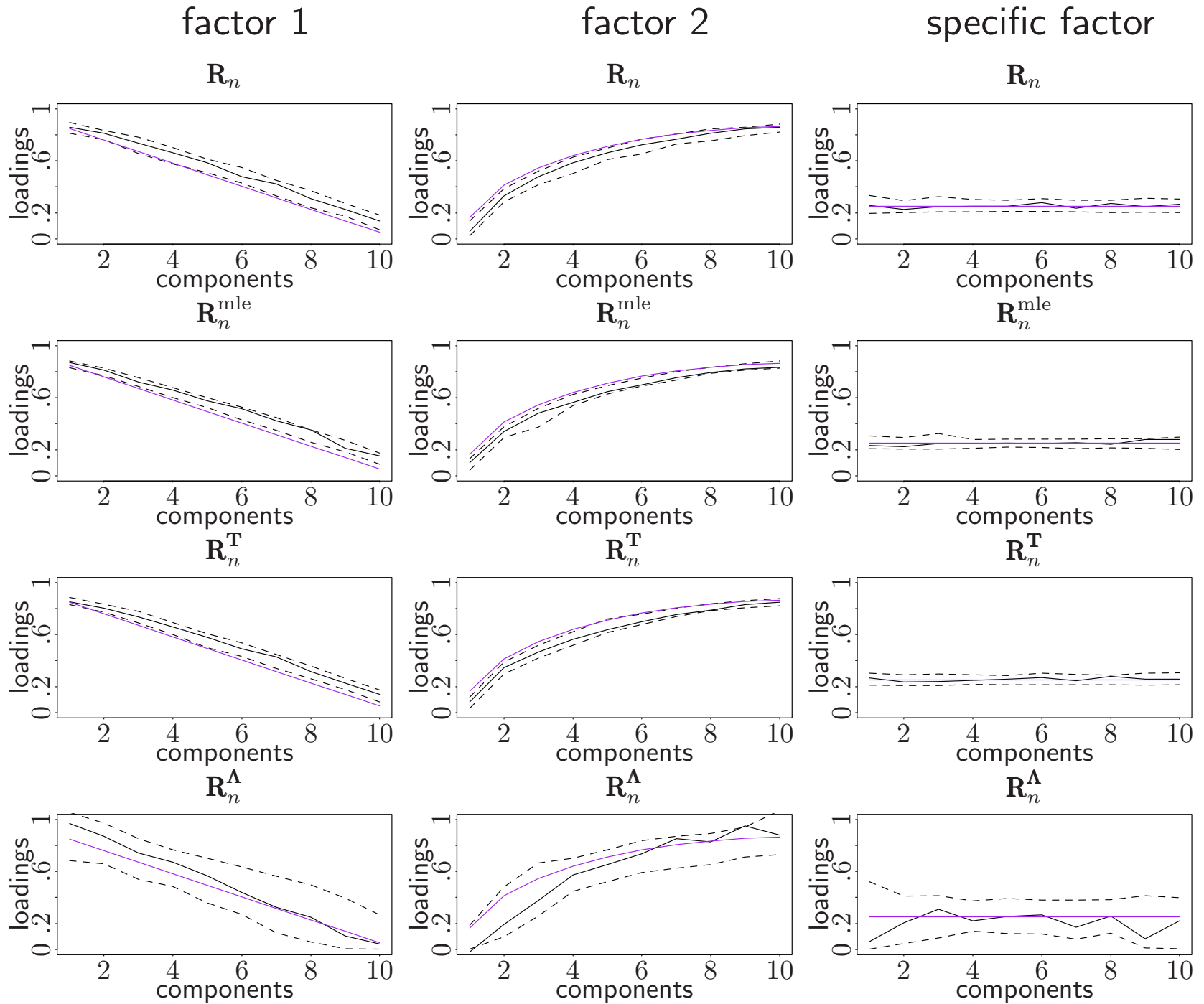
- Choose $d = 10$, $k = 2$ factors with loadings

component	1	2	3	4	5	6	7	8	9	10
$\mathbf{L}_{\cdot,1}$.85	.76	.67	.58	.49	.41	.32	.23	.14	.05
$\mathbf{L}_{\cdot,2}$.17	.41	.55	.64	.71	.77	.81	.84	.85	.86
$\text{diag}(\mathbf{V}^2)$.25	.25	.25	.25	.25	.25	.25	.25	.25	.25

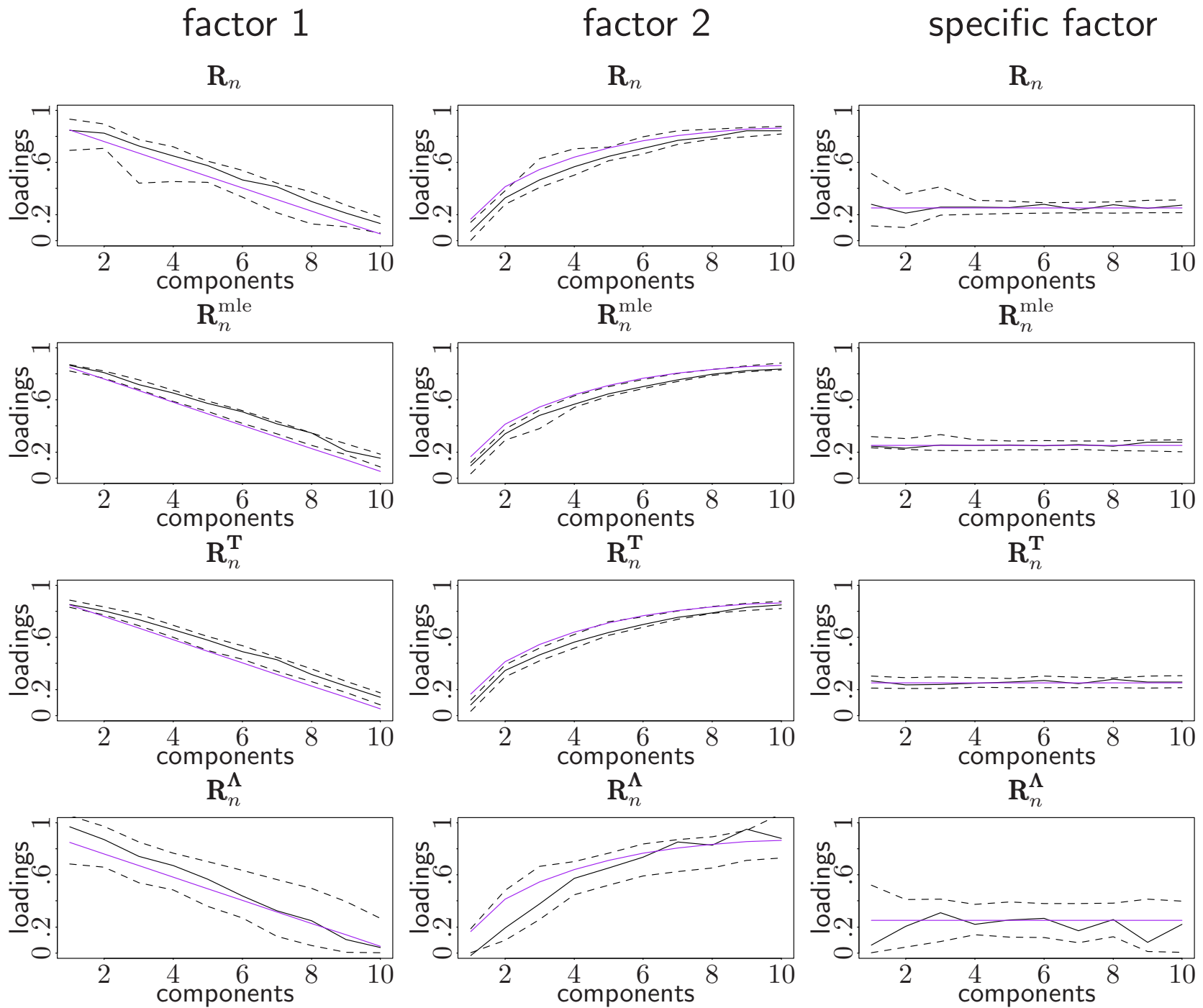
$(\Rightarrow \mathbf{L}\mathbf{L}^T + \mathbf{V}^2 = \mathbf{R}$ is a correlation matrix)

- consider factor model(s) (given before)
 1. $\mathbf{X} \stackrel{d}{=} \mathbf{L}\mathbf{f} + \mathbf{V}\mathbf{e} \sim t_d(\mathbf{0}, \mathbf{R}, \nu)$ with $\nu = 6$
 2. \mathbf{X} has same $t(\mathbf{0}, \mathbf{R}, \nu)$ dependence structure,
but different margins $F_i = t_{\nu_i}, \boldsymbol{\nu} = (3, \dots, 10)$
- Simulation length $n = 2000$, repeat 500 times
Plots of different estimation methods and 95%-CI's of loadings

$$\mathbf{X} \sim t_d(\mathbf{0}, \mathbf{R}, \nu)$$



X has t -dependence, different margins



Example

- Consider 8-dimensional set of data:
oil, s&p500, gbp, usd, chf, jpy, dkk and sek (exchange rates w.r.t. euro)
- Daily log-returns between May, 1985 to June, 2004 (n=4904)
- Apply factor analysis with different estimators as before

