

# **Model selection for FA/PCA**

Shikui Tu

**Department of Computer Science and  
Engineering, Shanghai Jiao Tong University**

**2021-04-25**

# Outline

- Recall
  - Probabilistic PCA, Factor Analysis (FA)
- Model selection for PCA/FA
- Practice on experimental comparisons

# PCA by minimizing MSE

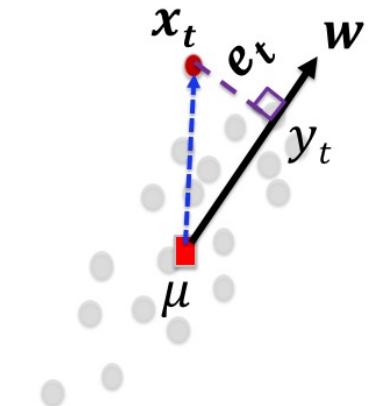
$$J(\mathbf{w}) = \frac{1}{N} \sum_{t=1}^N \|\mathbf{x}_t - (\mathbf{x}_t^T \mathbf{w}) \mathbf{w}\|^2$$

$$\begin{aligned} & \mathbf{x}_t^T \mathbf{x}_t - (\mathbf{x}_t^T \mathbf{w}) \mathbf{w}^T \mathbf{x}_t - \mathbf{x}_t^T (\mathbf{x}_t^T \mathbf{w}) \mathbf{w} + (\mathbf{x}_t^T \mathbf{w}) \mathbf{w}^T (\mathbf{x}_t^T \mathbf{w}) \mathbf{w} \\ &= \mathbf{x}_t^T \mathbf{x}_t - \mathbf{w}^T (\mathbf{x}_t \mathbf{x}_t^T) \mathbf{w} \end{aligned}$$

Introduce a Lagrange multiplier  $\lambda$

$$L(\{\mathbf{x}_t\}, \mathbf{w}) = J(\{\mathbf{x}_t\}, \mathbf{w}) - \lambda \cdot (\mathbf{w}^T \mathbf{w} - 1)$$

$$\frac{\partial J(\mathbf{w})}{\partial \mathbf{w}} - \lambda \cdot \frac{\partial (\mathbf{w}^T \mathbf{w} - 1)}{\partial \mathbf{w}} = -2(\Sigma_x \mathbf{w}) - \lambda \cdot 2\mathbf{w} = \mathbf{0}$$



$$\|\mathbf{w}\| = 1$$

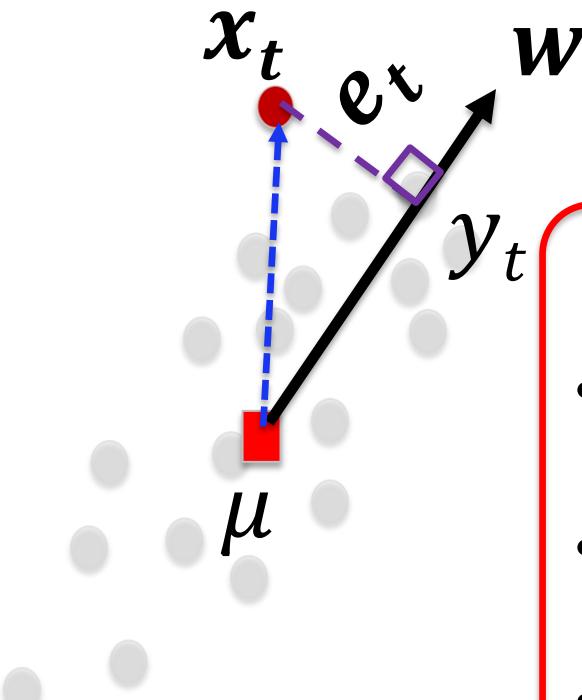
$$\Sigma_x = \frac{1}{N} \sum_{t=1}^N \mathbf{x}_t \mathbf{x}_t^T$$

$$\Sigma_x \mathbf{w} = (-\lambda) \cdot \mathbf{w}$$

Eigenvalues and Eigenvectors

# Factor Analysis (FA) model

Continuous latent variable  $y$



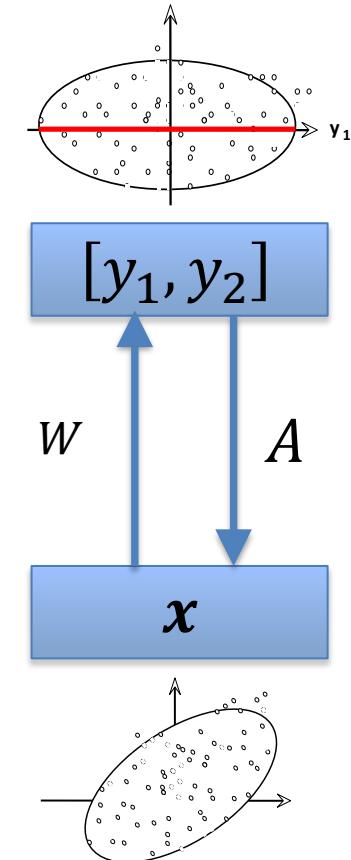
$$\|w\| = 1$$

$$y_t = x_t^T w$$

$$e_t = \|x_t - y_t w\|^2$$

For the  $t$ -th data point:

- Randomly sample a  $y_t$ :  
 $y_t \sim G(y|\mathbf{0}, \Sigma_y);$
- Randomly generate a noise  $e_t$   
 $e_t \sim G(e|0, \sigma^2 I)$
- Generate  $x_t$  by:  
$$x_t = Ay_t + \mu + e_t$$



# EM algorithm for FA

$$\text{E-Step: } p^{old}(\mathbf{y}|\mathbf{x}) = \frac{G(\mathbf{y}|0, I)G(\mathbf{x}|A\mathbf{y} + \boldsymbol{\mu}, \sigma^2 I)}{G(\mathbf{x}|\boldsymbol{\mu}, AA^T + \sigma^2 I)}$$

$$E[\mathbf{y}|\mathbf{x}] = W\mathbf{x} \quad W = A^T(AA^T + \sigma^2 I)^{-1}$$

$$E[\mathbf{y}\mathbf{y}^T|\mathbf{x}] = I - WA + W\mathbf{x}\mathbf{x}^TW^T$$

$$\text{M-Step: } \max Q(p^{old}(\mathbf{y}|\mathbf{x}), \Theta)$$

$$Q = \int p^{old}(\mathbf{y}|\mathbf{x}) \cdot \ln[G(\mathbf{y}|0, I)G(\mathbf{x}|A\mathbf{y} + \boldsymbol{\mu}, \sigma^2 I)] d\mathbf{y}$$

$$A^{new} = \left( \sum_{t=1}^N \mathbf{x}_t (E[\mathbf{y}|\mathbf{x}_t])^T \right) \left( \sum_{t=1}^N E[\mathbf{y}\mathbf{y}^T|\mathbf{x}_t] \right)^{-1}$$

$$\sigma^2{}^{new} = \frac{1}{Nd} Tr \left\{ \sum_{t=1}^N \{ \mathbf{x}_t \mathbf{x}_t^T - A^{new} E[\mathbf{y}|\mathbf{x}_t] \mathbf{x}_t^T \} \right\}$$

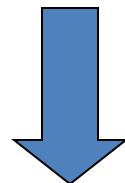
# Maximum likelihood FA implements PCA

$$p(\mathbf{y}) = G(\mathbf{y}|\mathbf{0}, I), \quad p(\mathbf{x}|\mathbf{y}) = G(\mathbf{x}|A\mathbf{y} + \boldsymbol{\mu}, \Sigma_e),$$

$$p(\mathbf{x}|\Theta) = \int p(\mathbf{y})p(\mathbf{x}|\mathbf{y})d\mathbf{y} = G(\mathbf{x}|\boldsymbol{\mu}, AA^T + \Sigma_e),$$

$$\max_{\Theta} \log \left\{ \prod_{t=1}^{N=1} p(\mathbf{x}_t | \Theta) \right\}$$

Maximum Likelihood



$$\Sigma_e = \sigma_e^2 \mathbf{I}_n$$

assume  $\boldsymbol{\mu} = \mathbf{0}$

PCA

$$\begin{cases} \hat{\mathbf{A}}_{n \times m}^{ML} = \mathbf{U}_{n \times m} (\mathbf{D}_m - \hat{\sigma}_e^2)^{\frac{1}{2}} \mathbf{R}^T, & \mathbf{D}_m = \text{diag}[s_1, \dots, s_m] \\ \hat{\sigma}_e^{2,ML} = \frac{1}{n-m} \sum_{i=m+1}^n s_i, \end{cases}$$

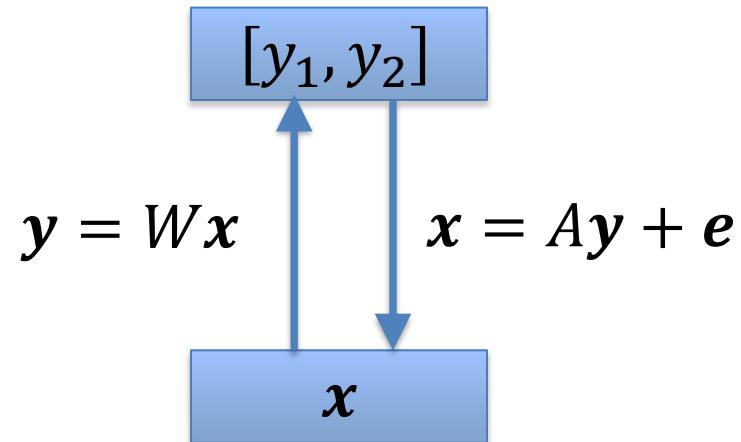
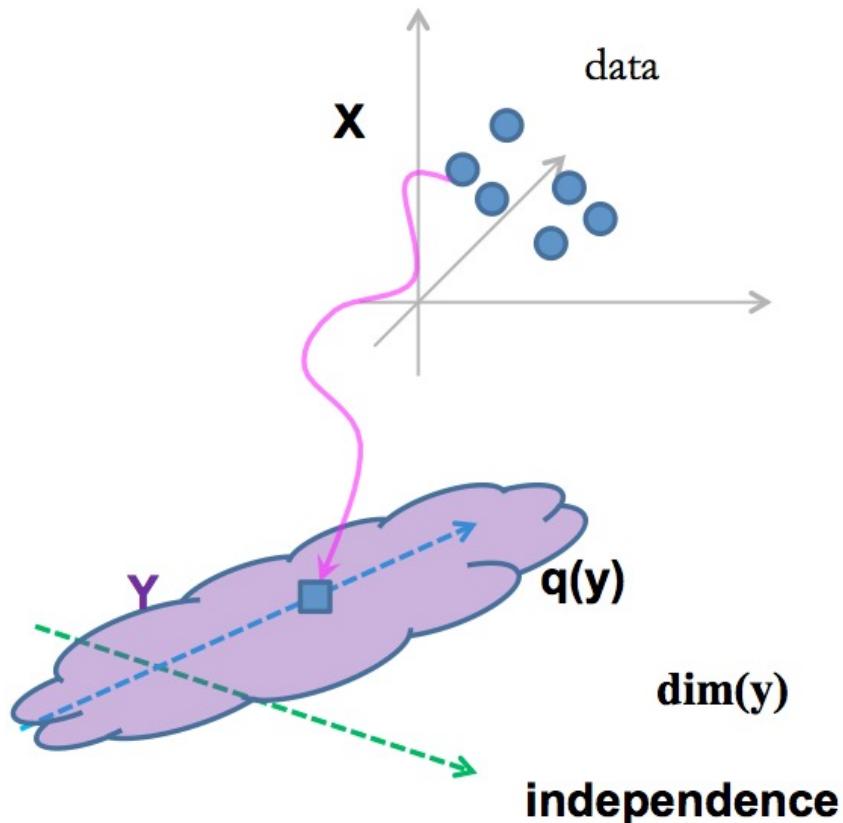
$\mathbf{U}$  is eigenvectors of sample cov.

# Outline

- Recall
  - Probabilistic PCA, Factor Analysis (FA)
- **Model selection for PCA/FA**
- Practice on experimental comparisons

# Dimensionality reduction

$$m = \dim(y) = ?$$



$$n = \dim(x) > m$$

Determining  $m$  is a model selection problem.

# Choosing $m$ by keeping 99% variance

$$m = \dim(\mathbf{y}) = ?$$

- Pick  $m$  to be the smallest value so that 99% of variance is retained, i.e.,

$$\frac{\frac{1}{N} \sum_{t=1}^N \|x_t - \hat{x}_t\|^2}{\frac{1}{N} \sum_{t=1}^N \|x_t\|^2} \leq 0.01 \quad (1\%)$$

Average squared projection error

Total variation in the data

Where  $\hat{x}_t$  is a reconstruction of  $x_t$  by

$$\hat{x}_t = W y_t = W(W^T x_t) \quad \hat{x}_t = y_t w = (x_t^T w)w$$

# Choosing $m$ by keeping 99% variance

- Eigen-decomposition

$$m = \dim(\mathbf{y}) = ?$$

$$\Sigma_x \approx U\Lambda U^T$$

$$U = [\mathbf{u}_1, \dots, \mathbf{u}_m]^T$$

$$\Lambda = \text{diag}[\lambda_1, \dots, \lambda_m]$$

$$\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_m$$

Covariance matrix (assume  $\mathbf{x}_t$  zero mean):

$$\Sigma_x = \frac{1}{N} \sum_{t=1}^N \mathbf{x}_t \mathbf{x}_t^T = \frac{1}{N} \mathbf{X} \mathbf{X}^T$$

$$\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_N]_{n \times N}$$

- Equivalently, we pick

Sum of selected eigenvalues

$$\frac{\sum_{i=1}^m \lambda_i}{\sum_{j=1}^n \lambda_j} \leq 0.01 \quad (1\%)$$

Sum of all eigenvalues

# Question

- Why are they equivalent?

$$\frac{\frac{1}{N} \sum_{t=1}^N \|x_t - \hat{x}_t\|^2}{\frac{1}{N} \sum_{t=1}^N \|x_t\|^2} \leq 0.01$$



$$\frac{\sum_{i=1}^m \lambda_i}{\sum_{j=1}^n \lambda_j} \leq 0.01$$

# The effects of $\dim(\mathbf{y})$

Assume  $\mu = \mathbf{0}$

$$\mathbf{x} = \mathbf{a}_1 y_1 + \mathbf{a}_2 y_2 + \cdots + \mathbf{a}_{m^*} y_{m^*} + \mathbf{e}$$

$$\mathbf{x} = \mathbf{A}\mathbf{y} + \mathbf{e}$$

$$\mathbf{y} = [y_1, \dots, y_{m^*}]_{m^* \times 1}^T$$

$$p(\mathbf{y}) = G(\mathbf{y} | \mathbf{0}, I),$$

Under-fitting (big bias):  $m < m^*$

Fitting error

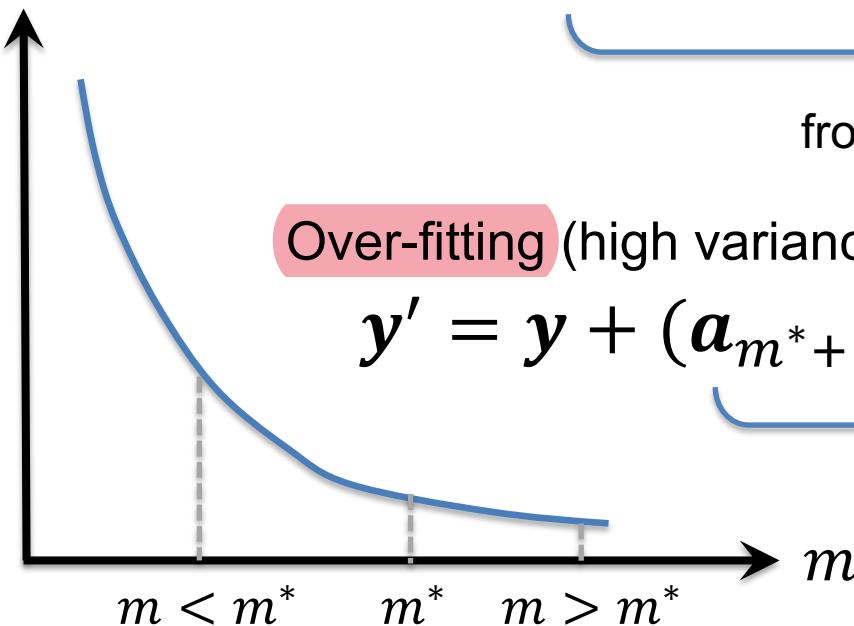
$$\mathbf{e}' = \mathbf{a}_{m+1} y_{m+1} + \cdots + \mathbf{a}_{m^*} y_{m^*} + \mathbf{e}$$

from signal  $\mathbf{y}$

Over-fitting (high variance):  $m > m^*$

$$\mathbf{y}' = \mathbf{y} + (\mathbf{a}_{m^*+1} y_{m^*+1} + \cdots + \mathbf{a}_m y_m)$$

from noise  $\mathbf{e}$



# Bias-variance decomposition

We want to find a function  $\hat{f}(y)$ , that approximates the true function  $f(y)$  as well as possible:

$$x = f(y) + \epsilon$$

$$E[\epsilon] = 0$$

$$E[x] = f$$

The expected error

$$Var[x] = Var[\epsilon]$$

$$\begin{aligned} E \left[ (x - \hat{f}(y))^2 \right] &= E[x^2 + \hat{f}^2 - 2x\hat{f}] \\ &= Var[x] + (E[x])^2 + Var[\hat{f}] + (E[\hat{f}])^2 - 2fE[\hat{f}] \\ &= Var[x] + Var[\hat{f}] + (f - E[\hat{f}])^2 \\ &= Var[\epsilon] + \text{Var}[\hat{f}] + \text{Bias}[\hat{f}]^2 \end{aligned}$$



# Model selection for FA

Probabilistic model

$$p(X_N | \Theta_K)$$

Candidate models:

$$\Theta_1 \subseteq \Theta_2 \subseteq \dots \subseteq \Theta_K \subseteq \dots$$

## Factor Analysis

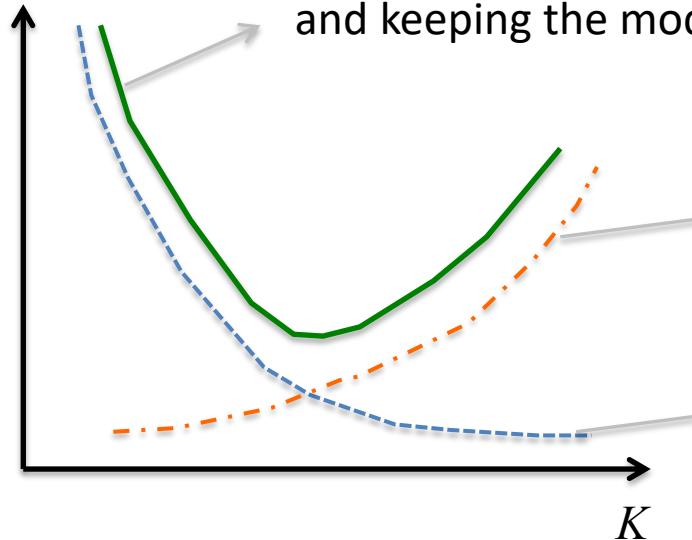
$$p(\mathbf{y}) = G(\mathbf{y}|\mathbf{0}, I),$$

$$p(\mathbf{x}|\mathbf{y}) = G(\mathbf{x}|A\mathbf{y} + \boldsymbol{\mu}, \Sigma_e),$$

$$\begin{aligned} p(\mathbf{x}|\Theta) &= \int p(\mathbf{y})p(\mathbf{x}|\mathbf{y})d\mathbf{y} \\ &= G(\mathbf{x}|\boldsymbol{\mu}, AA^T + \Sigma_e), \end{aligned}$$

$$K = m = \dim(\mathbf{y})$$

Criterion



A trade-off between fitting the data well and keeping the model simple

Models become more and more complex as  $K$  increases.

Negative log-likelihood (or fitting error) to make sure the model fit the data well.

Akaike's Information Criterion (AIC)

$$\ln p(X_N | \hat{\Theta}_K) - d_k$$

Bayesian Information Criterion (BIC)

$$\ln p(X_N | \hat{\Theta}_K) - \frac{1}{2} d_k \ln N$$

$d_k$ : number of free parameters  
 $N$ : sample size

# Two-stage model selection

Given a sample set  $\mathcal{X}_N = \{\mathbf{x}_t\}_{t=1}^N$ , the task of FA modeling consists of parameter learning and model selection, i.e., selecting the number of factors  $m$ . Given  $m$ , the parameter set  $\Theta_m$  is usually estimated by Maximum Likelihood (ML) learning. Including it as a subtask, model selection problem is usually tackled by the following classical two-stage procedure:

- **Stage I:** Compute  $\hat{\Theta}_m = \hat{\Theta}(\mathcal{X}_N, m)$  for each  $m \in [m_{low}, m_{up}]$  with  $m_{low}$  and  $m_{up}$  given. Normally,  $\hat{\Theta}_m$  is the ML estimator,

$$\begin{aligned}\hat{\Theta}_m^{ML} &= \arg \max_{\Theta_m} \ln p(\mathcal{X}_N | \Theta_m) \\ &= \arg \min_{\Theta_m} \mathcal{E}_L(\mathcal{X}_N | \Theta_m),\end{aligned}\quad (2)$$

where  $\mathcal{E}_L(\mathcal{X}_N | \Theta_m) = -\frac{2}{N} \ln p(\mathcal{X}_N | \Theta_m)$  is denoted as **NLL** (negative log-likelihood).

- **Stage II:** Estimate  $\hat{m} = \arg \min_m \mathcal{E}_{Cri}$ , where  $\mathcal{E}_{Cri}$  is formulated according to an information criterion (Cri), e.g.,

$$\mathcal{E}_{Cri}(\mathcal{X}_N, \hat{\Theta}_m) = \mathcal{E}_L(\mathcal{X}_N, \hat{\Theta}_m) + \frac{\rho_N d_m}{N} \quad (3)$$

$$d_m = nm + 1 - \frac{m(m-1)}{2} \quad (4)$$

$$\rho_N = \begin{cases} 0; & \text{for NLL} \\ 2; & \text{for AIC} \\ \ln N; & \text{for BIC/MDL} \\ \ln N + 1; & \text{for CAIC} \\ 2 \ln(\ln N); & \text{for HQC} \end{cases} \quad (5)$$

where  $d_m$  by eq.(4) is the number of free parameters of FA model over the real number field, and  $d_m = m(2n - m) + 1$  when considering complex number field. In the

# Variational FA

$$\Theta_m = \{\mathbf{A}, \sigma_e^2\}$$

$$\mathbf{x} = \mathbf{Ay} + \boldsymbol{\mu} + \mathbf{e},$$

$$a_k^\alpha = b_k^\alpha = a^\phi = b^\phi = 10^{-3}$$

$$\boldsymbol{\mu} = \mathbf{0}, \Sigma_e = \sigma_e^2 \mathbf{I}_n$$

$$p(\mathbf{x}|\mathbf{y}) = G(\mathbf{x}|\mathbf{Ay} + \boldsymbol{\mu}, \Sigma_e),$$

$$p(\mathbf{y}) = G(\mathbf{y}|\mathbf{0}, \Sigma_y)$$

$$p(\mathbf{A}|\boldsymbol{\alpha}) = \prod_{k=1}^m G(\mathbf{a}_k|\mathbf{0}, \frac{1}{\alpha_k^\alpha} \mathbf{I}_n)$$

$$p(\boldsymbol{\alpha}|\boldsymbol{a}^\alpha, \boldsymbol{b}^\alpha) = \prod_{k=1}^m \Gamma(\alpha_k|a_k^\alpha, b_k^\alpha)$$

$$p(\varphi) = \Gamma(\varphi|a^\varphi, b^\varphi), \varphi = (\sigma_e^2)^{-1}$$

$$\mathcal{F}(q_Y, q_{\mathbf{A}}, q_{\boldsymbol{\alpha}}, q_\varphi) = \int q_Y q_{\mathbf{A}} q_{\boldsymbol{\alpha}} q_\varphi \ln \frac{p(\mathcal{X}_N, Y, \Theta | m)}{q_Y q_{\mathbf{A}} q_{\boldsymbol{\alpha}} q_\varphi} dY d\mathbf{A} d\boldsymbol{\alpha} d\varphi dY$$

Stage I: get  $\mathcal{F}(m)$ ,  $m \in [m_{low}, m_{up}]$ ;

Stage II: estimate  $\hat{m} = \arg \max_m \mathcal{F}(m)$ .

# VBEM for FA

**Initialization:** get  $\{q_{\mathbf{A}}, q_{\boldsymbol{\alpha}}, q_{\varphi}\}$  by random initialization.

**E-Step:** get  $q_Y$  based on  $q_{\theta} = q_{\mathbf{A}} q_{\boldsymbol{\alpha}} q_{\varphi}$ :

$$q_Y = \prod_t G(\mathbf{y}_t | \Sigma_y \langle \mathbf{A}^T \boldsymbol{\varphi} \rangle_{q_{\theta}} \mathbf{x}_t, \Sigma_y),$$

where  $\Sigma_y = (\mathbf{I} + \langle \mathbf{A}^T \boldsymbol{\varphi} \mathbf{A} \rangle_{q_{\theta}})^{-1}$ ,  $\langle \bullet \rangle_{q_{\theta}} = E_{q_{\theta}}[\bullet]$ .

**M-Step:** get  $q_{\theta_i}$  based  $q_{\theta_\ell}$  and  $q_Y$ , where

$$q_{\theta_i}, q_{\theta_\ell} \in \{q_{\mathbf{A}}, q_{\boldsymbol{\alpha}}, q_{\varphi}\}, \ell \neq i :$$

$$q_{\mathbf{A}} = \prod_{j=1}^n G(\mathbf{a}_j | \mu_{a,j}, \Sigma_{a,j}),$$

$$\Sigma_{a,j} = \left( \langle \boldsymbol{\varphi} \rangle_{q_{\varphi}} \sum_t \langle \mathbf{y}_t \mathbf{y}_t^T \rangle_{q_Y} + \text{diag}(\langle \boldsymbol{\alpha} \rangle_{q_{\boldsymbol{\alpha}}}) \right)^{-1},$$

$$\mu_{a,j} = \Sigma_{a,j} \langle \boldsymbol{\varphi}_j \rangle_{q_{\varphi}} \left( \sum_t \mathbf{x}_{jt} \langle \mathbf{y}_t \rangle_{q_Y} \right),$$

$$q_{\boldsymbol{\alpha}} = \prod_{k=1}^m \Gamma \left( \alpha_k | a_k^\alpha + \frac{n}{2}, b_k^\alpha + \frac{\langle \mathbf{a}_k \mathbf{a}_k^T \rangle_{q_{\mathbf{A}}}}{2} \right),$$

$$q_{\varphi} = \Gamma \left( \varphi | a^\varphi + \frac{N}{2}, b^\varphi + \frac{\sum_t \langle (\mathbf{x}_{jt} - \mathbf{a}_j^T \mathbf{y}_t)^2 \rangle_{q_{\mathbf{A}} q_Y}}{2} \right).$$

# Outline

- Recall
  - Probabilistic PCA, Factor Analysis (FA)
- Model selection for PCA/FA
- **Practice on experimental comparisons**

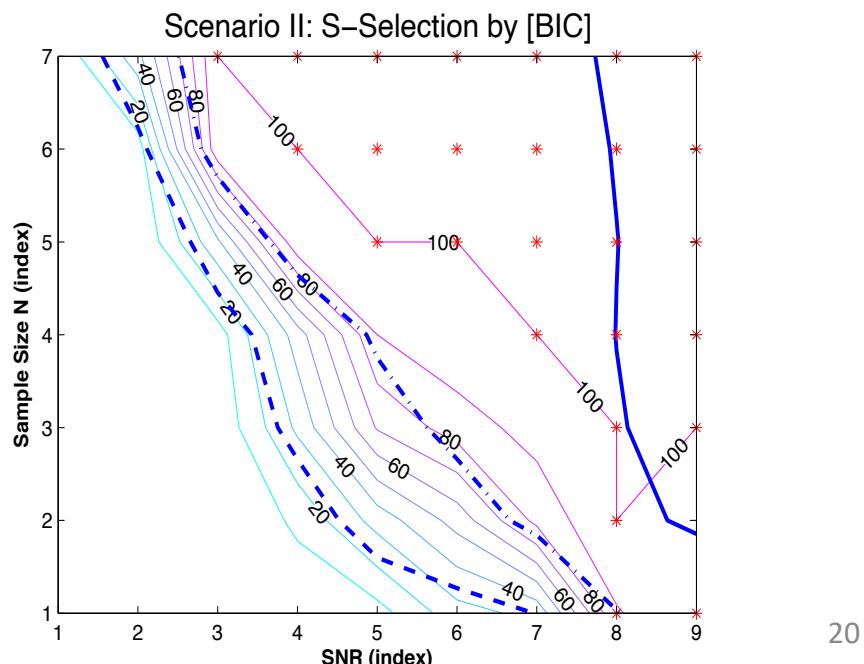
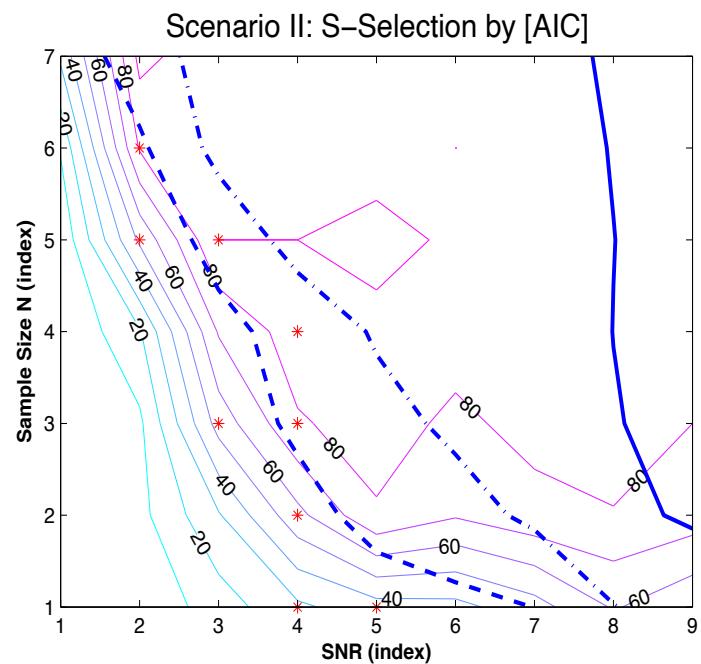
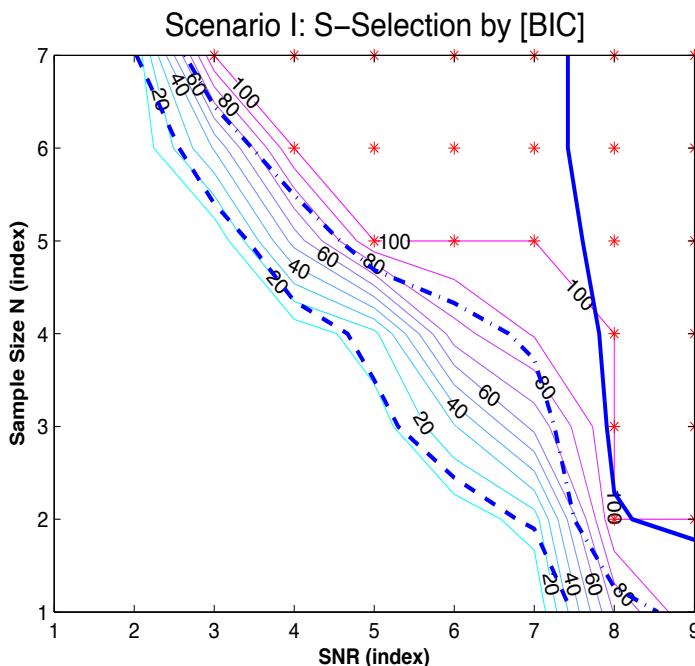
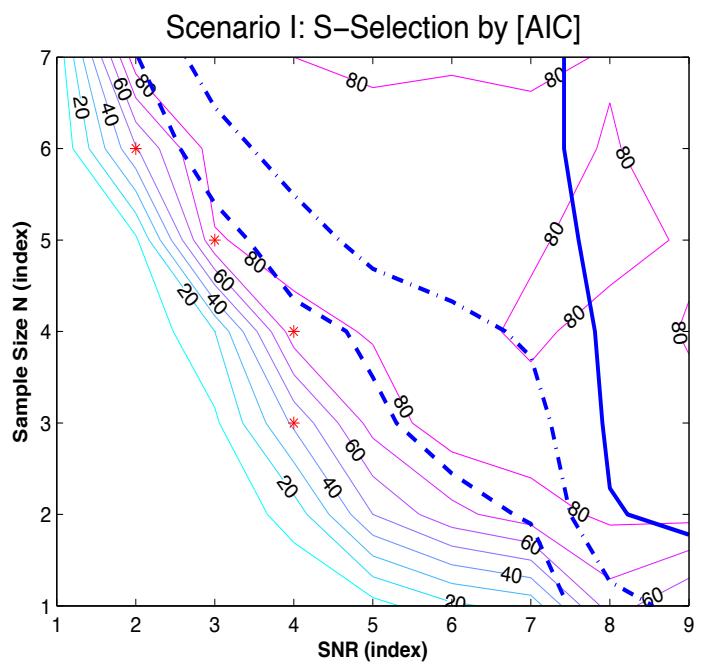
# An experimental study

## Experimental settings

id	features $f$	values $V(f)$
$f_1$	sample size $N$	$\{25, 50, 75, 100, 200, 400, 800\}$
$f_2$	SNR $\gamma_o = \lambda_{m^*} / \sigma_e^2 + 1$	$\{1.2, 1.5, 2, 2.5, 3,$ $3.5, 4, 8, 16\}$
$f_3$	$\lambda_i \sim \mathcal{U}([\lambda_{min}, \lambda_{max}])$	$[1, 1], [1, 10], [1, 50]$
$f_4$	$\{n, m\} = \{\dim(x), \dim(y)\}$	$\{15, 5\}$

Each setting:  $\Omega(f_1, f_2, f_3, f_4)$ ,  $f_4 : \{n = 15, m = 5\}$ .

scenarios ( $\tau$ )	criteria: {AIC,BIC, ...}	categories(c)
<b>I</b> $(f_3 : \mathcal{U}([1, 1]))$	$\forall c$ riterion,	<b>U:</b> $\hat{m} < m^*$
<b>II</b> $(f_3 : \mathcal{U}([1, 10]))$	$\forall (f_1, f_2) \in V(f_1) \times V(f_2)$	<b>S:</b> $\hat{m} = m^*$
<b>III</b> $(f_3 : \mathcal{U}([1, 50]))$	$R_{c,\tau}^{(cri)}(f_1, f_2)$	<b>O:</b> $\hat{m} > m^*$



# Thank you!