

# Decomposition methods for explorative neuroimaging

Lars Kai Hansen

DTU Informatics  
Technical University of Denmark

Co-workers:

Morten Mørup, Kristoffer Madsen, Finn Å. Nielsen, Mads Dyrholm, Stephen Strother, Rasmus Olsson, Thomas Kolenda, Niels Mørch, Ulrik Kjems, Sidse Arnfred, Benny Lautrup.



*Do not multiply causes!*

## OUTLINE

- Hypothesis testing vs knowledge discovery
  - Generalizability
- Factor models - Linear hidden variable representations
- Principal component analysis (PCA)
  - Understanding the limits to learning in high-dimensional data
  - Heuristics to heal overfitting in poor SNR's: Re-scaling projections
- Independent component analysis (ICA)
  - Statistical modeling and generalization
  - Generalizations: Convolutional mixing, shift
- More generalizations
  - NMF
  - Multiway modeling

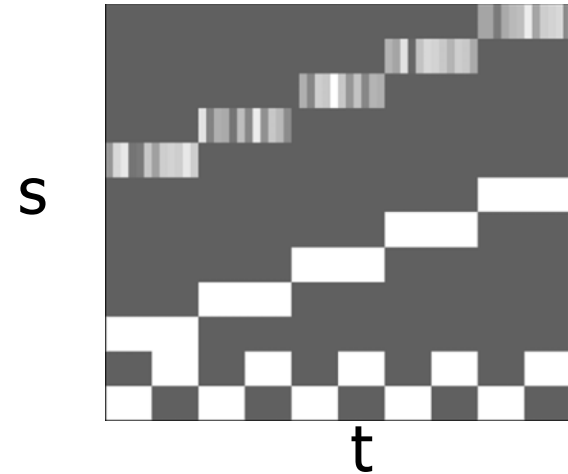


# Multivariate neuroimaging models

Neuroimaging aims at extracting the mutual information between stimulus and response.

- Stimulus: Macroscopic variables, "design matrix" ...  $s(t)$
- Response: Micro/meso-scopic variables, the neuroimage ...  $x(t)$
- Mutual information is stored in the joint distribution ...  $p(x,s)$ .

*Often  $s(t)$  is assumed known....unsupervised methods consider  $s(t)$  or parts of  $s(t)$  "hidden".....*



# Multivariate neuroimaging models

- Univariate models (SPM, fMRI time series models etc)

$$p(x, s) = p(x | s)p(s) = \prod_j p(x_j | s) \cdot p(s)$$



- Multivariate models (PCA, PLS, ICA, ANN etc)

$$p(x, s) = p(s | x)p(x)$$

- Modeling from data ( $D$ ) w. parameterized function families

$$p(s | x) \sim p(s | x, D) \sim p(s | x, \hat{\theta}), \quad \hat{\theta} = \hat{\theta}(D)$$

$$p(x) \sim p(x | D) \sim p(x | \hat{\theta}),$$

# Early multivariate neuroimage analyses

- **Principal component analysis**

Regional analysis (Moeller and Strother, *JCBFM* 1991), Spatial (Friston et al., *JCBFM* 1993), Selection of subspace dimension (Hansen et al. *NeuroImage* 1999)

- **Linear multivariate models**

Worsley et al. (*NeuroImage* 1997), Discriminants (Mørch et al. *IPMI* 1997)

- **Non-linear multivariate models**

Artif. neural networks (Lautrup et al., 1995; Mørch et al., *IPMI* 1997)

Independent components (McKeown et al., *PNAS* 1998)

Spatio-temporal clustering (Ding et al. ISMR 1994, Toft et al, *HBM Conf. 1997*, Goutte et al., *NeuroImage* 1999) Meta analysis (Goutte et al. *HBM* 2001)

- **Plurality and similarity**

ROC curves for multiple models (Lange et al., *NeuroImage* 1999)

Consensus of activation maps (Hansen et al., *NeuroImage* 2001)

# Generative model

$$\mathbf{x} = \mathbf{A}\mathbf{s} + \boldsymbol{\varepsilon}, \quad \boldsymbol{\varepsilon} \sim N(\mathbf{0}, \boldsymbol{\Sigma})$$

$$p(\mathbf{x} | \mathbf{A}, \boldsymbol{\theta}) = \int p(\mathbf{x} | \mathbf{A}, \mathbf{s}, \boldsymbol{\Sigma}) p(\mathbf{s} | \boldsymbol{\theta}) d\mathbf{s}$$

$$p(\mathbf{x} | \mathbf{A}, \mathbf{s}, \boldsymbol{\Sigma}) = |2\pi\boldsymbol{\Sigma}|^{-1/2} e^{-\frac{1}{2}(\mathbf{x}-\mathbf{A}\mathbf{s})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x}-\mathbf{A}\mathbf{s})}$$

Source distribution:

PCA: ... normal

ICA: ... other

IFA: ... Gauss. Mixt.

$$\text{PCA: } \boldsymbol{\Sigma} = \sigma^2 \cdot \mathbf{1},$$

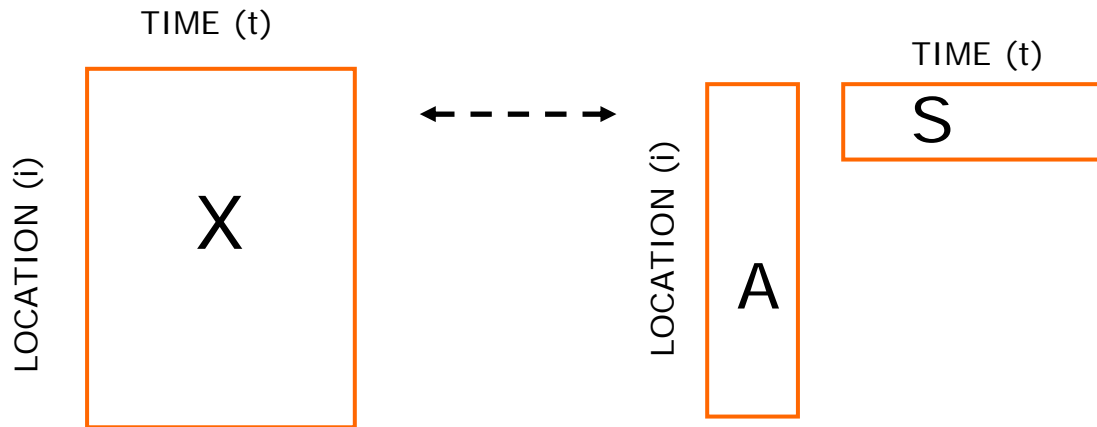
$$\text{FA: } \boldsymbol{\Sigma} = \mathbf{D}$$

S known:	GLM
(1-A <sup>-1</sup> ) sparse:	SEM
S,A positive:	NMF

Højten-Sørensen, Winther, Hansen,  
Neural Comp (2002), Neurocomputing (2002)

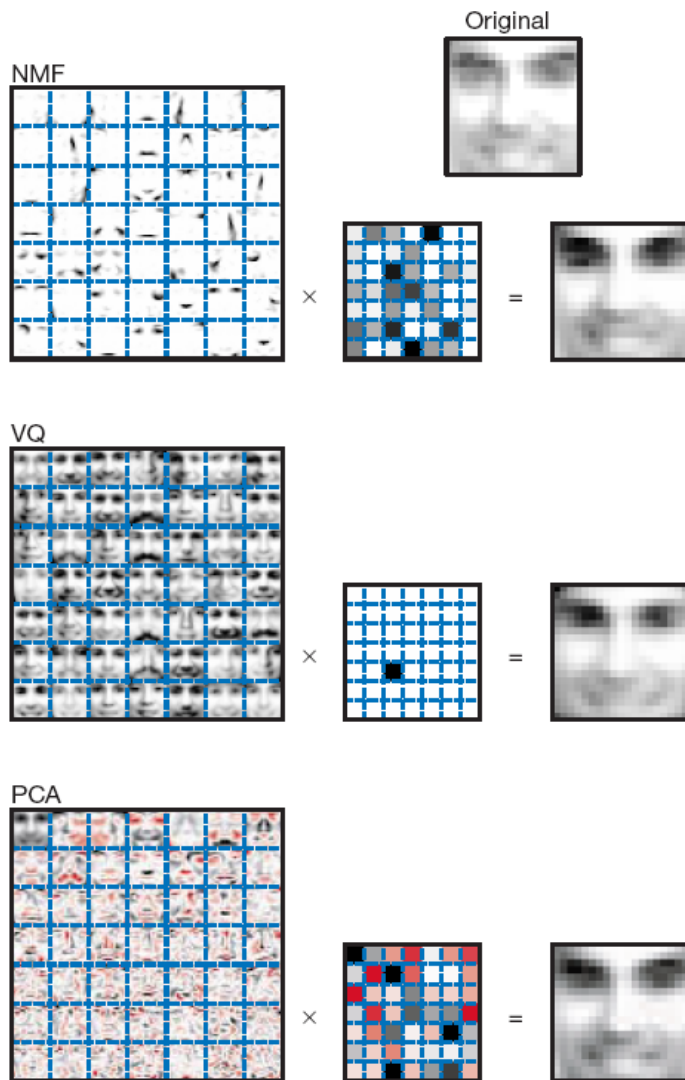
# Factor models

- Represent a datamatrix by a low-dimensional approximation



$$X(i, t) \approx \sum_{k=1}^K A(i, k) S(k, t)$$

# Matrix factorization: SVD/PCA, NMF, Clustering



**Figure 1** Non-negative matrix factorization (NMF) learns a parts-based representation of faces, whereas vector quantization (VQ) and principal components analysis (PCA) learn holistic representations. The three learning methods were applied to a database of  $m = 2,429$  facial images, each consisting of  $n = 19 \times 19$  pixels, and constituting an  $n \times m$  matrix  $V$ . All three find approximate factorizations of the form  $V \approx WH$ , but with three different types of constraints on  $W$  and  $H$ , as described more fully in the main text and methods. As shown in the  $7 \times 7$  montages, each method has learned a set of  $r = 49$  basis images. Positive values are illustrated with black pixels and negative values with red pixels. A particular instance of a face, shown at top right, is approximately represented by a linear superposition of basis images. The coefficients of the linear superposition are shown next to each montage, in a  $7 \times 7$  grid, and the resulting superpositions are shown on the other side of the equality sign. Unlike VQ and PCA, NMF learns to represent faces with a set of basis images resembling parts of faces.

## Learning the parts of objects by non-negative matrix factorization

Daniel D. Lee\* & H. Sebastian Seung\*†

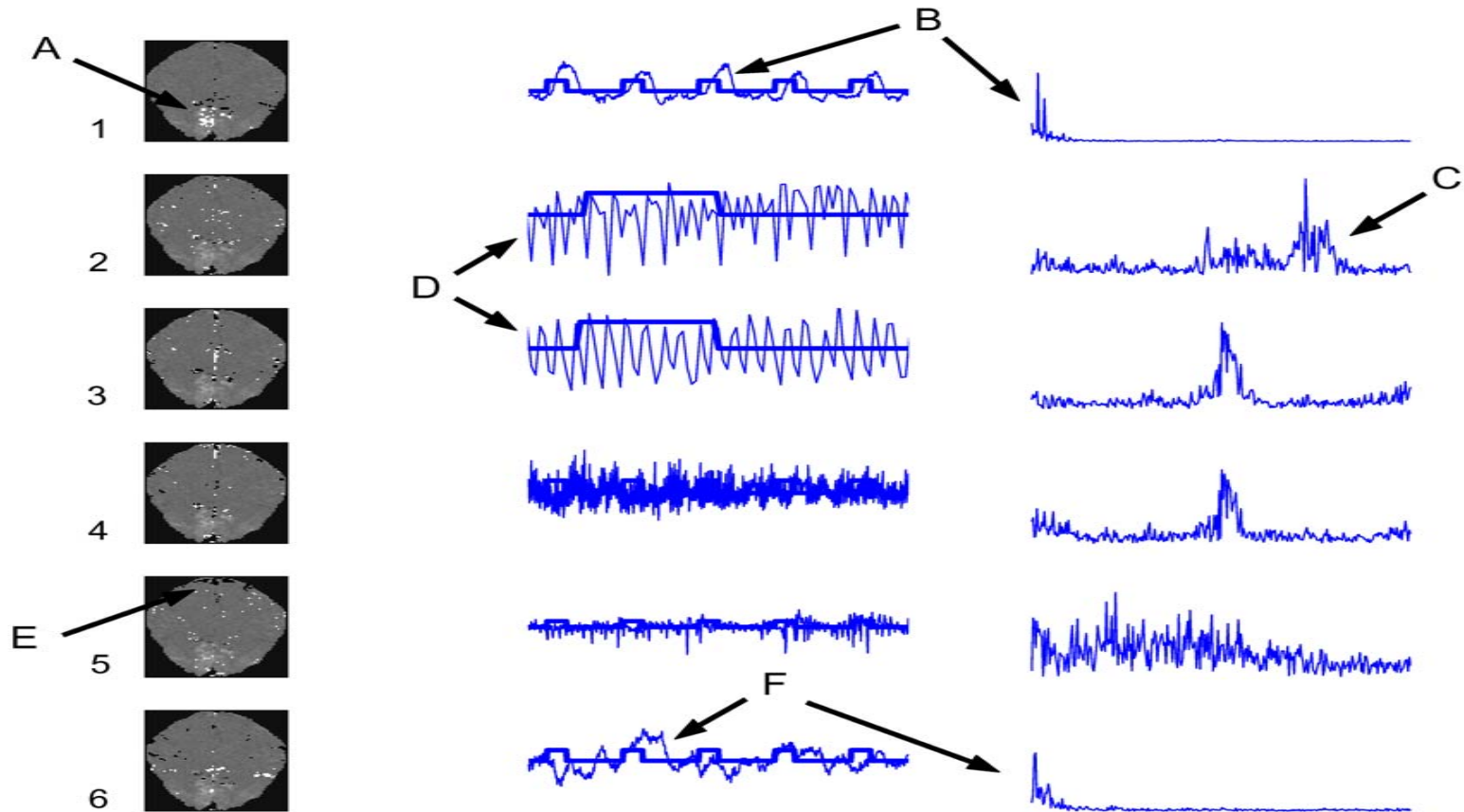
\* Bell Laboratories, Lucent Technologies, Murray Hill, New Jersey 07974, USA

† Department of Brain and Cognitive Sciences, Massachusetts Institute of Technology, Cambridge, Massachusetts 02139, USA

NATURE | VOL 401 | 21 OCTOBER 1999 | www.nature.com



# ICA: Assume $S(k,t)$ , $S(k',t)$ statistically independent



(McKeown, Hansen, Sejnowski, Curr. Op. in Neurobiology (2003))

# Generalizability

*Do not multiply causes!*

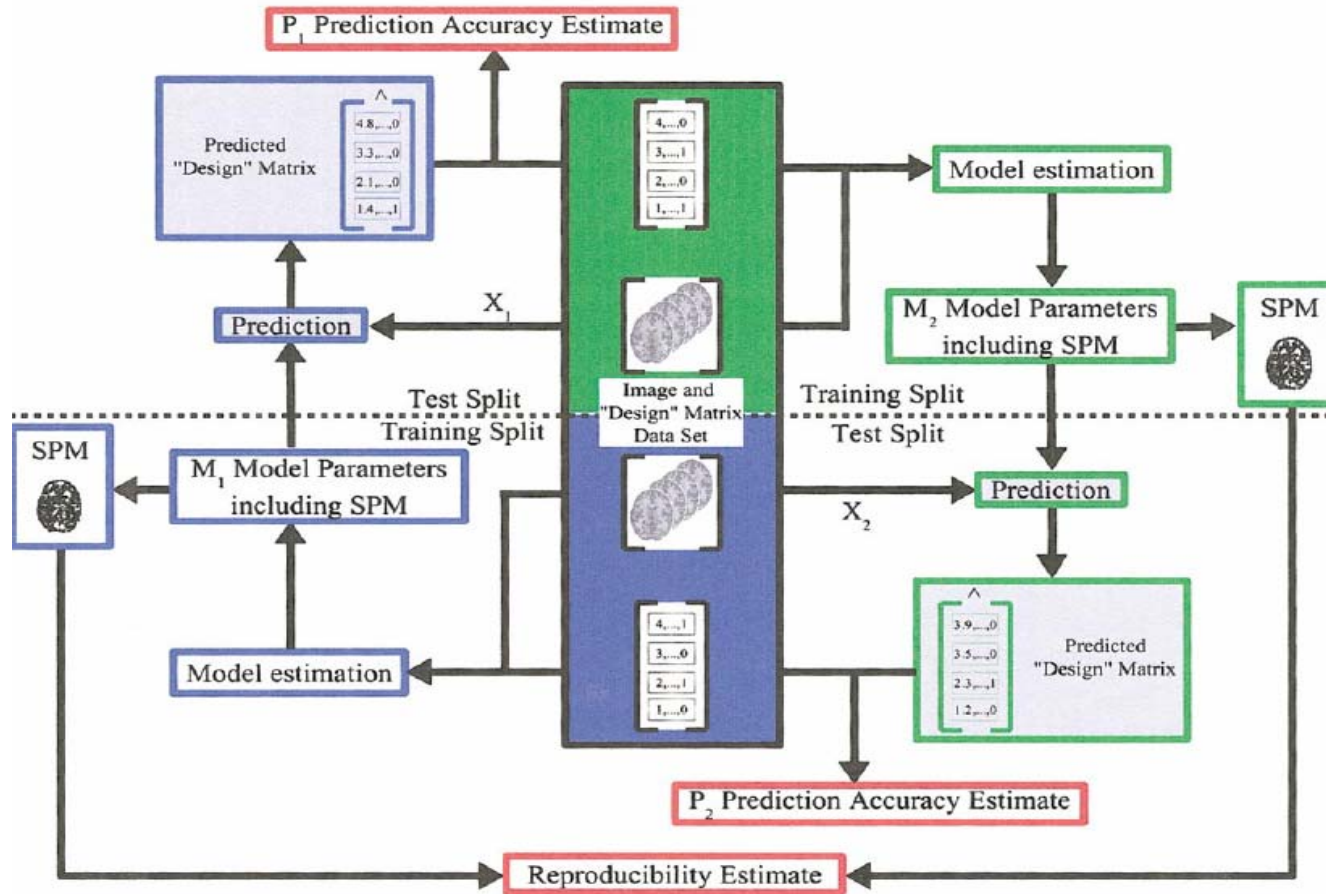


- Generalizability is defined as *the expected performance on a random new sample* ... the mean performance of a model on a "fresh" data set is an unbiased estimate of generalization
- Typical loss functions:

$$\langle -\log p(\mathbf{s} | \mathbf{x}, D) \rangle, \quad \langle -\log p(\mathbf{x} | D) \rangle,$$
$$\langle (\mathbf{s} - \hat{\mathbf{s}}(D))^2 \rangle, \quad \left\langle \log \frac{p(\mathbf{s}, \mathbf{x} | D)}{p(\mathbf{s} | D)p(\mathbf{x} | D)} \right\rangle$$

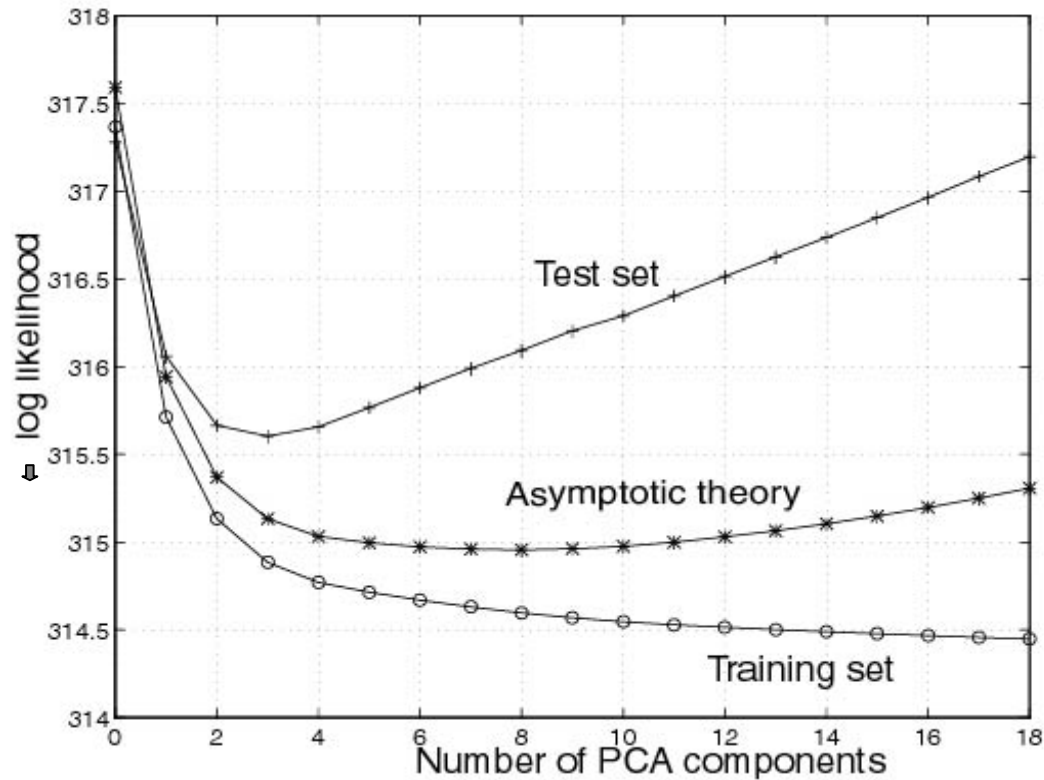
- Note: Even for hidden variable models we can "predict"!
- Results can be presented as "bias-variance trade-off curves" or "learning curves"

# NPAIRS: Reproducibility of parameters



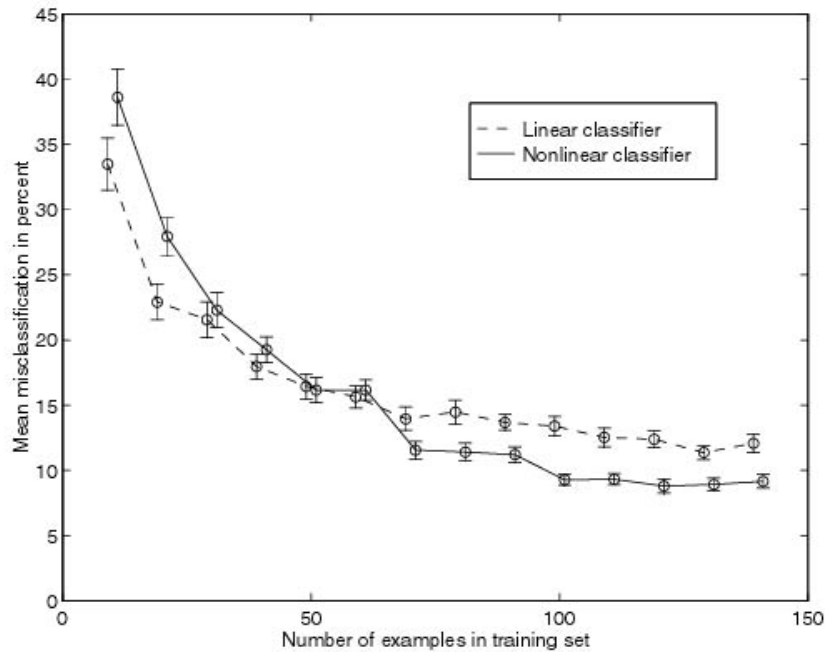
NeuroImage: Hansen et al (1999), Hansen et al (2000), Strother et al (2002), Kjems et al. (2002), LaConte et al (2003), Strother et al (2004)

# Bias-variance trade-off as function of PCA dimension

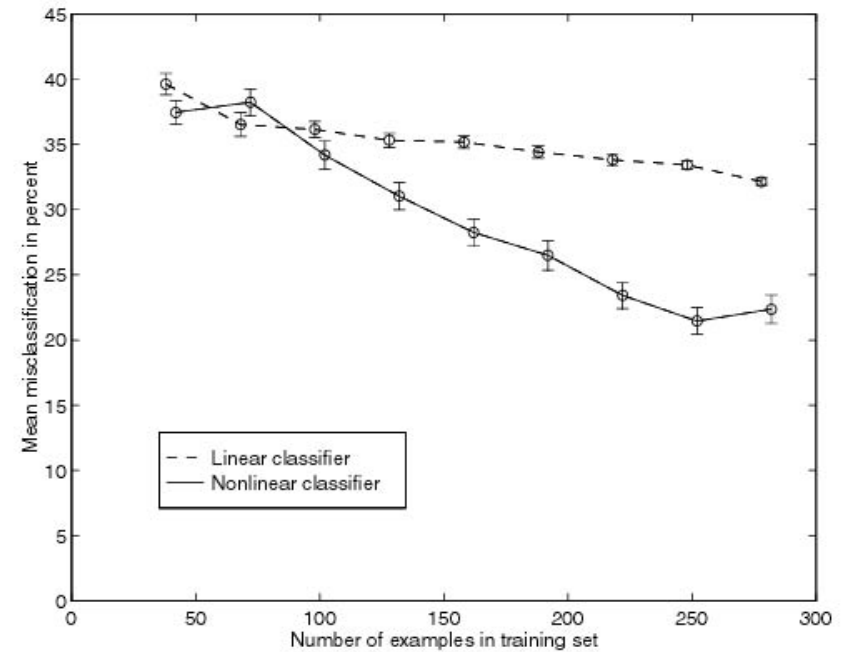


Hansen et al. *NeuroImage* (1999)

# Learning curves multivariate models: Within group generalizability



PET



fMRI

Finger tapping, analysed by PCA dimensional reduction and Fisher LD / ML Perceptron. Mørch et al. *IPMI* (1997)... "first mind reading in fMRI"

# Modeling the generalizability of SVD

- Rich physics literature on "retarded" learning

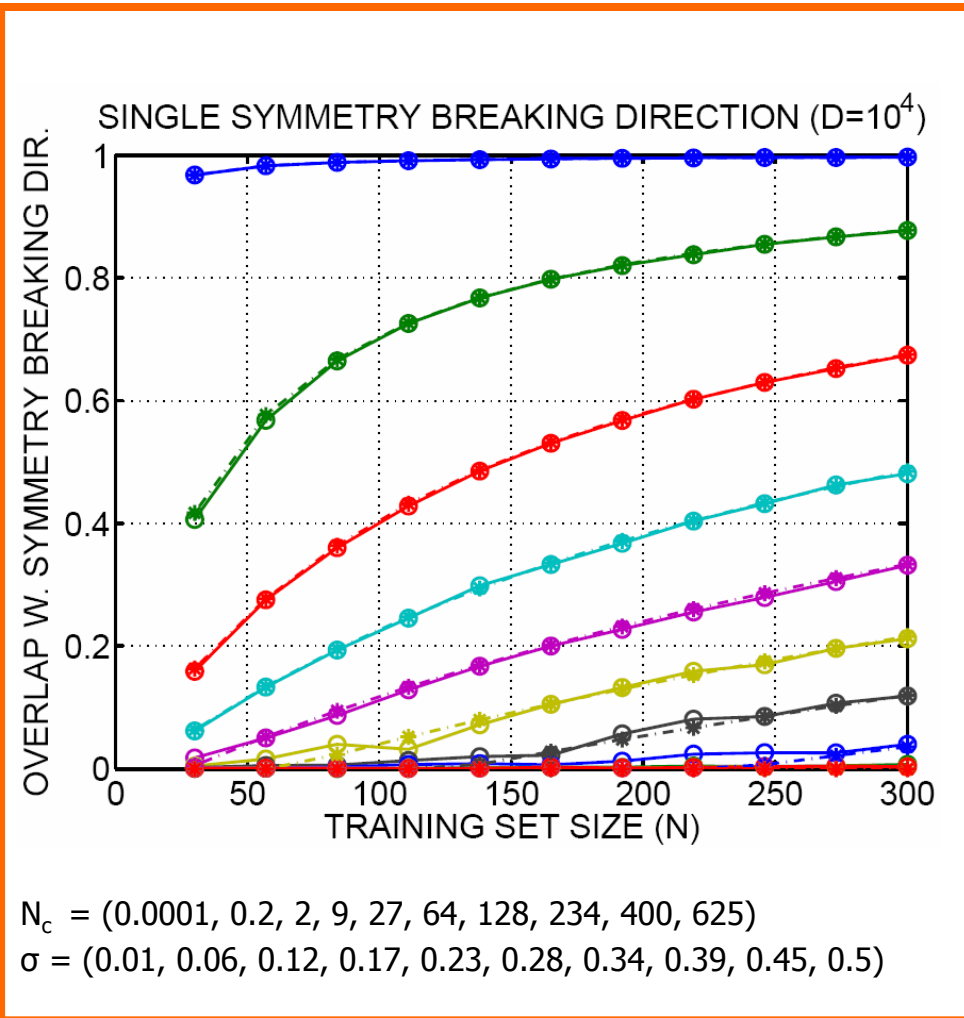
- Universality**

- Generalization for a "single symmetry breaking direction" is a function of ratio of  $N/D$  and signal to noise  $S$
- For subspace models-- a bit more complicated -- depends on the component SNR's and eigenvalue separation
- For a single direction, the mean squared overlap  $R^2 = \langle (u_1^T u_0)^2 \rangle$  is computed for  $N, D \rightarrow \infty$

$$R^2 = \begin{cases} (\alpha S^2 - 1) / S(1 + \alpha S) & \alpha > 1/S^2 \\ 0 & \alpha \leq 1/S^2 \end{cases}$$

$$\alpha = N/D \quad S = 1/\sigma^2 \quad N_c = D/S^2$$

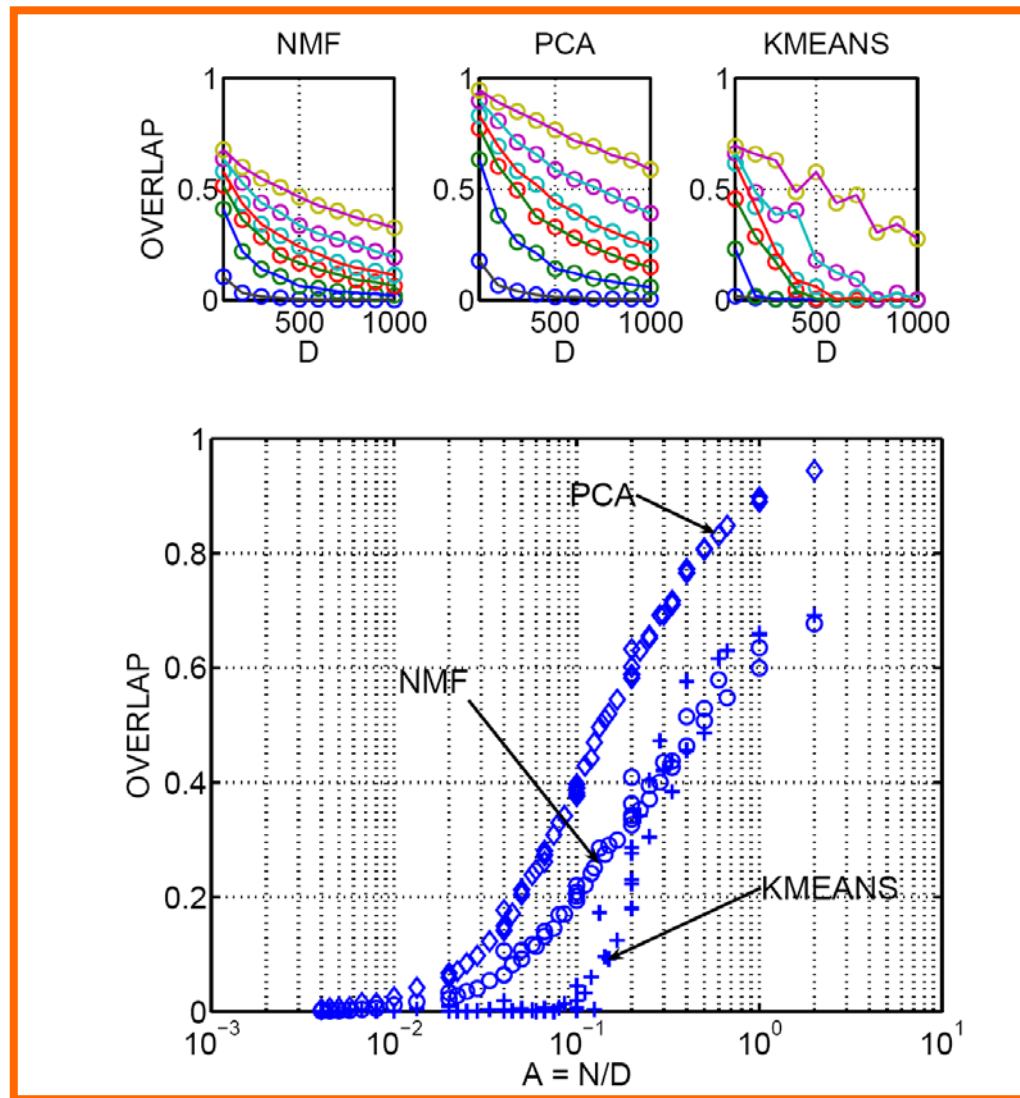
Hoyle, Rattray: Phys Rev E **75** 016101 (2007)



# Universality in PCA, NMF, Kmeans

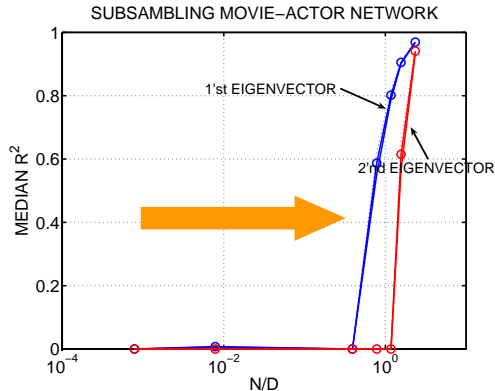
- Looking for universality by simulation
  - learning two clusters in white noise.
- Train  $K=2$  component factor models.
- Measure overlap between line of sight and plane spanned by the two factors.

Experiment  
Variable:  $N, D$   
Fixed: SNR

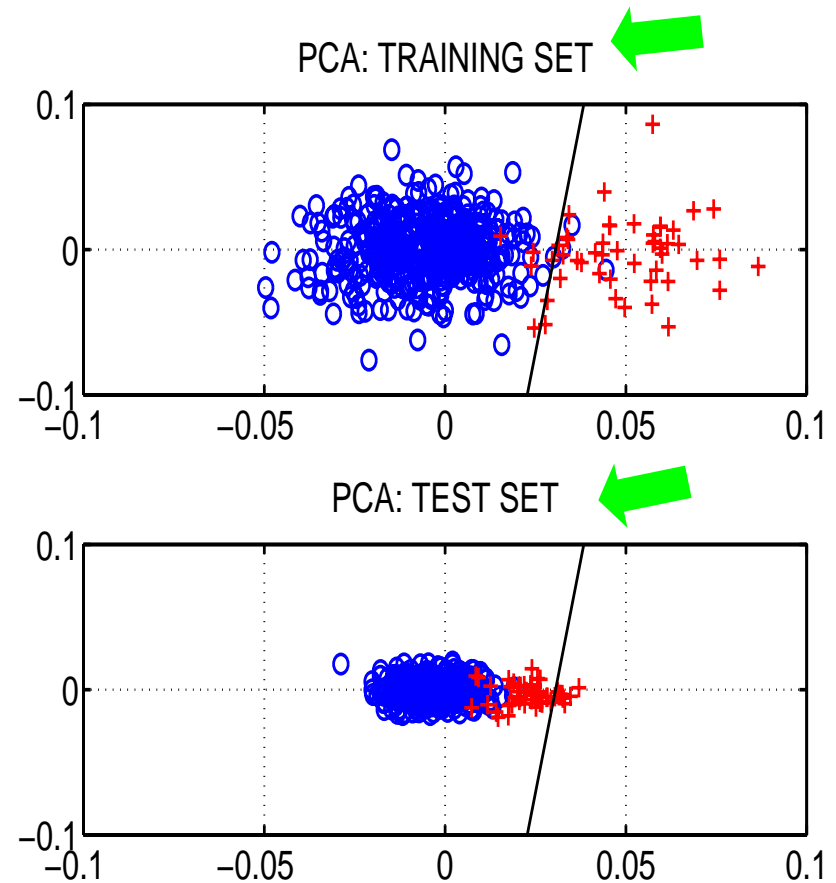


# Restoring the generalizability of SVD

- Now what happens if you are on the slope of generalization, i.e.,  $N/D$  is just beyond the transition to retarded learning ?

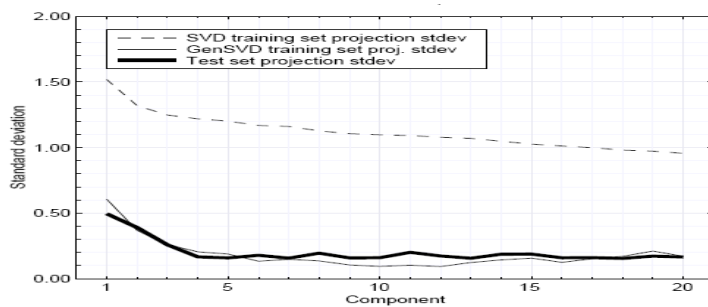
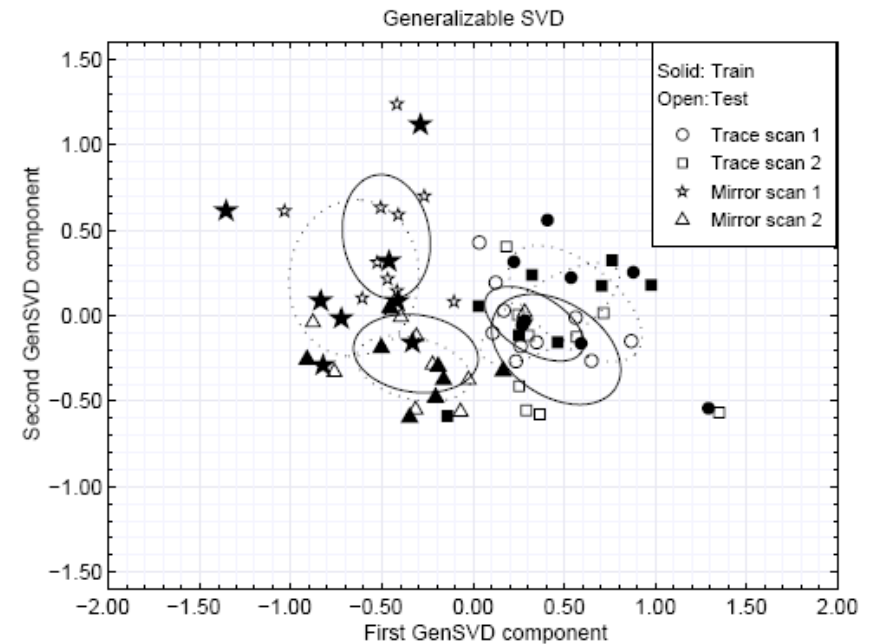
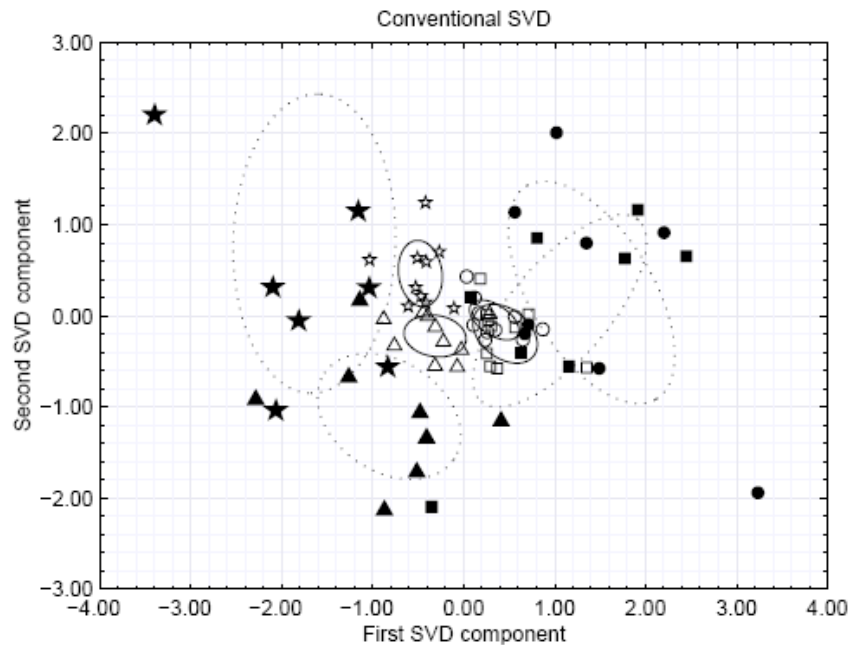


- The estimated projection is offset, hence, future projections will be too small!
- ...problem if discriminant is optimized for unbalanced classes in the training data!





# Heuristic: Leave-one-out re-scaling of SVD test projections

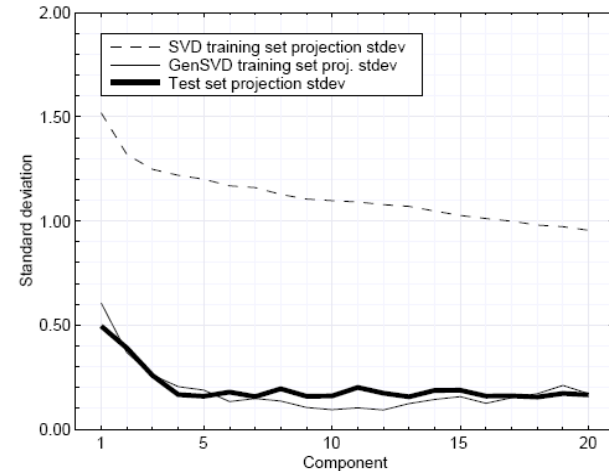


$N=72, D=2.5 \cdot 10^4$

Kjems, Hansen, Strother: "Generalizable SVD for Ill-posed data sets" NIPS (2001)

# Re-scaling the component variances

- Possible to compute the new scales by leave-one-out doing  $N$  SVD's of size  $N \ll D$



Compute  $\mathbf{U}_0 \mathbf{\Lambda}_0 \mathbf{V}_0^T = \text{svd}(X)$  and  $\mathbf{Q}_0 = [\mathbf{q}_j] = \mathbf{\Lambda}_0 \mathbf{V}_0^T$   
 foreach  $j = 1 \dots N$

$$\bar{\mathbf{q}}_{-j} = \frac{1}{N-1} \sum_{j' \neq j} \mathbf{q}_{j'}$$

$$\text{Compute } \mathbf{B}_{-j} \mathbf{\Lambda}_{-j} \mathbf{V}_{-j}^T = \text{svd}(\mathbf{Q}_{-j} - \bar{\mathbf{Q}}_{-j})$$

$$\mathbf{z}_j = \mathbf{B}_{-j} \mathbf{B}_{-j}^T (\mathbf{q}_j - \bar{\mathbf{q}}_{-j})$$

$$\hat{\lambda}_i^2 = \frac{1}{N-1} \sum_j z_{ij}^2$$

Kjems, Hansen, Strother: NIPS (2001)

# Identifiability problem in SVD/PCA

Linear model + normality = failure!

$$x(j, t) = \sum_k A(j, k) s(k, t)$$

$$\begin{aligned} R &= \frac{1}{T} \sum_t x(j, t) x(j', t) = \langle xx' \rangle = A \langle ss' \rangle A' = AA' \\ &= AUU' A' = AU(AU)' = BB' \end{aligned}$$

If columns of  $A$  are in general position,  
 $A$  can not be recovered by PCA  
Call ICA: the rotation buster



# Factor analysis with S "independent" - ICA issues in fMRI analysis

- Basic assumptions:
  - Spatial or temporal independence?
  - Which ICA model, higher order stat or temporal corr to bust rotation?
  - Are confounds independent of design variables?
- Statistical issues
  - How many components to use?
  - Systematic evaluation in terms of generalization

Consensus and model averaging

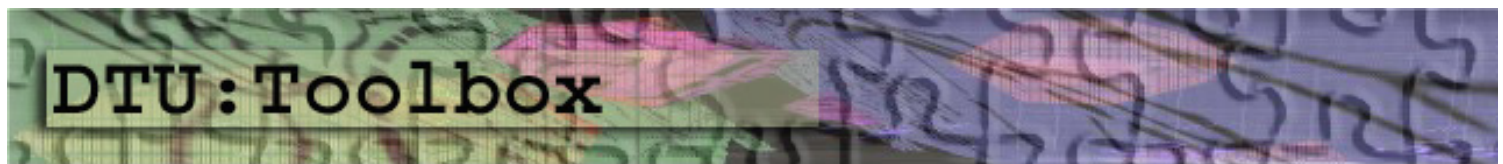
(1000+ papers/abstracts on ICA+fMRI)

# Independent Component Analyses (select. refs.)

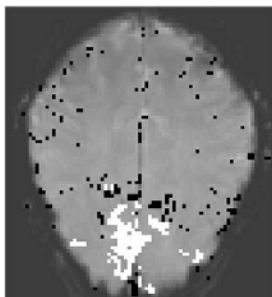
- First non-linear approach (Herault&Jutten, 1985)
- Blind signal separation (Cardoso, Comon, 1987-)
  
- Temporal decorrelation (Molgedey&Schuster, 1994)
- Infomax: High impact work by (Bell&Sejnowski, 1995)
- Infomax applications in EEG: Scott Makeig et al. (1996) "EEG-Lab"
  
- Application for fMRI expl. analysis (McKeown et al, 1998)
- High-dimensional mixtures (Hansen&Larsen, 1998)
- Noisy image mixtures => estimated sources  
are non-linear in the measurement! (Hansen, 2000)
- Bayesian ICA (Højen-Sørensen, Hansen & Winther 2001, 2002a, 2002b)
  
- Reviews of ICA and fMRI (e.g. McKeown, Hansen, Sejnowski, 2003)
- Special issue of IEEE Trans Bio. Engineering (Calhoun & Adali, 2006)

# DTU:ICA toolbox ([www.imm.dtu.dk/cisp](http://www.imm.dtu.dk/cisp))

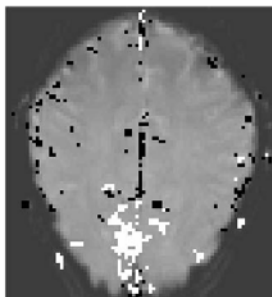
- **Infomax/Maximum likelihood**  
Bell & Sejnowski (1995), McKeown et al (1998)
- **Dynamic Components**  
Molgedey-Schuster (1994), Petersen et al (2001)
- **Mean Field ICA**  
Højen-Sørensen et al. (2001,2002)
  
- **Features:**
  - v Number of components (BIC)
  - v Parameter tuning
  - v Binary and mixing constraints (A)
  - v Demo scripts incl. fMRI data



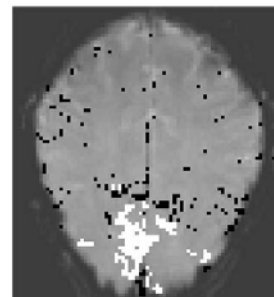
## Spatial mode ICA



Attias

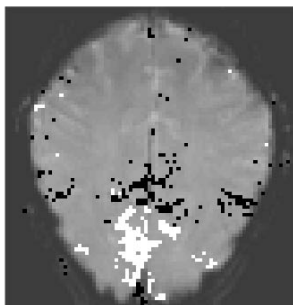


Decorr

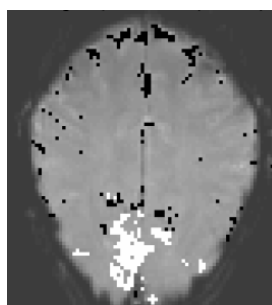


Max Lik.

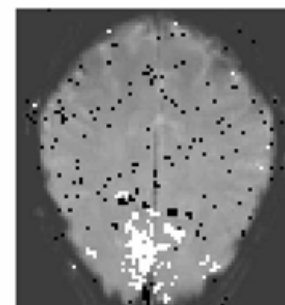
## Temporal mode ICA



Attias

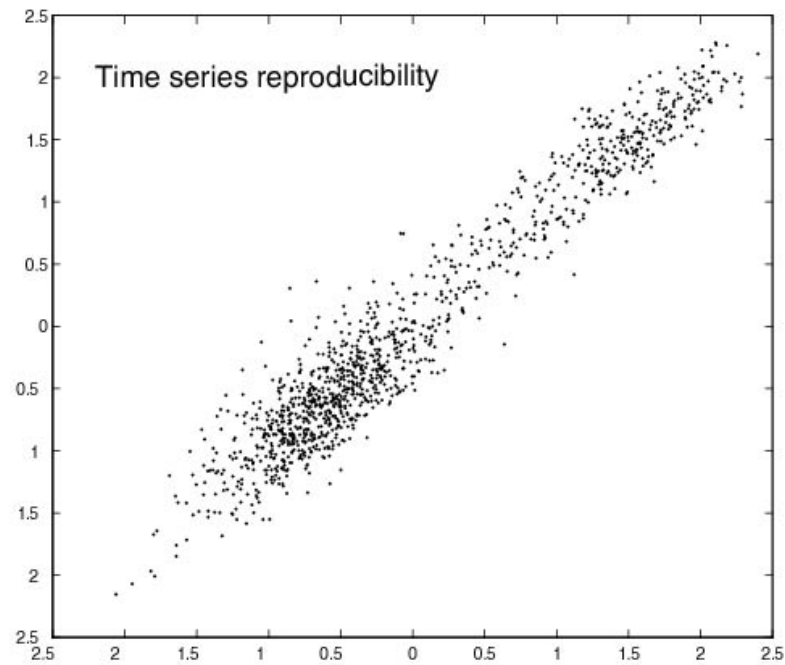


Decorr



Max Lik.

# Reproducibility of the time series recovered from spatial vs temporal ICA



Petersen et al. ICA2000





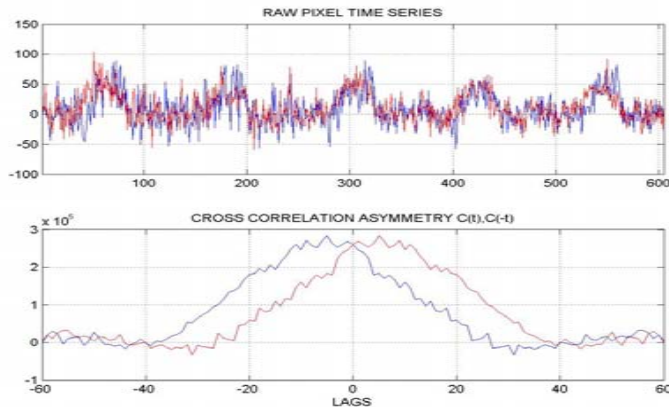
“Mr. Osborne, may I be excused? My brain is full.”

# Challenges for the linear factor model

- Temporal structure in networks -> Convolutional ICA
- Group study, repeat trials -> Multiway methods
- ~~Causal models –sparse ICA~~

# fMRI: Delayed activation in visual cortex

## Two voxel temporal cross-correlation



$$\langle x(j,t)x(j',t+\tau) \rangle = \sum_{k,k'} A(j,k)A(j',k') \langle s(k,t)s(k',t+\tau) \rangle$$

$$\langle x(j,t)x(j',t+\tau) \rangle = \sum_k A(j,k)A(j',k)C_k(\tau)$$

$$\Sigma_{j,j'}^{xx}(\tau) = \Sigma_{j,j'}^{xx}(-\tau)$$

## Recovery of source signals from an unknown, linear convolutive mixture

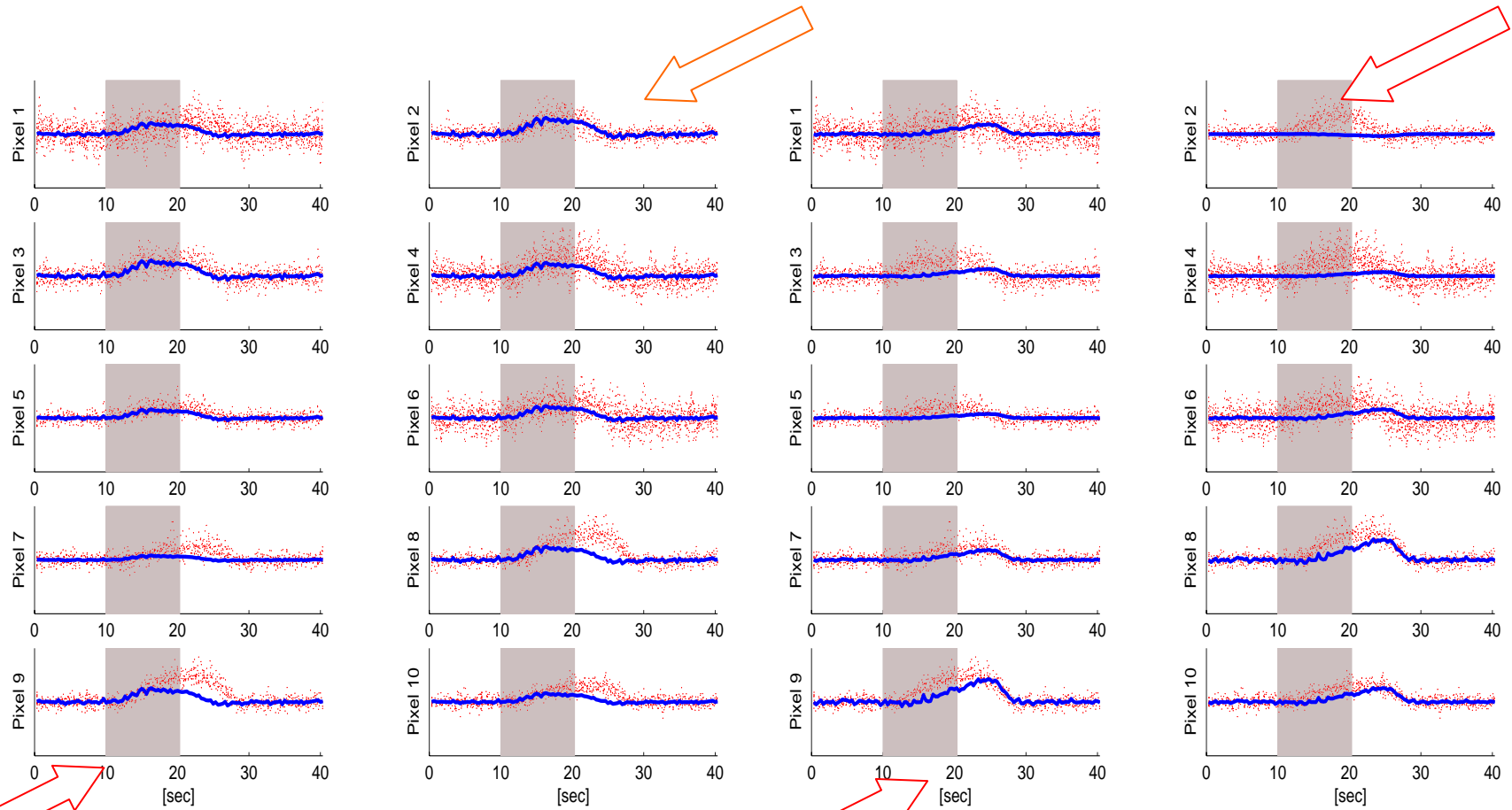
$$x(j,t) = \sum_{k,\tau} A(j,k,\tau) s(k,t-\tau)$$

Hansen & Dyrholm: Prediction matrix (Proc MLSP, 2003)

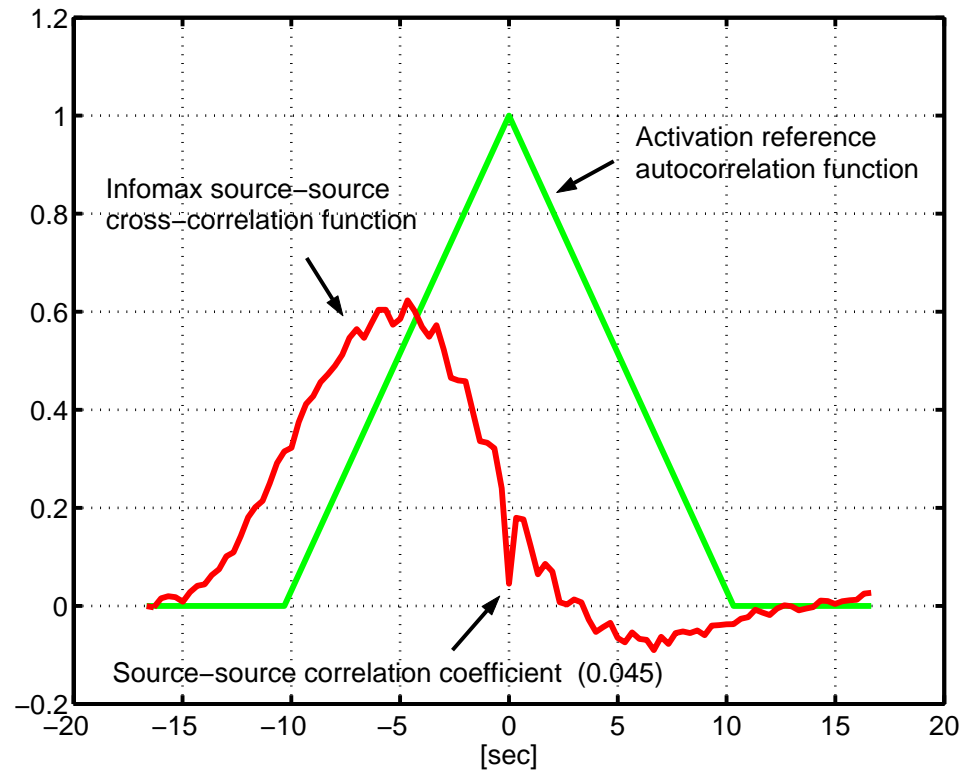
Olsson & Hansen: Kalman (Proc NIPS 2004), J. Mach. Learning Res. (2006)

Dyrholm & Hansen: "CICAAR" (Proc ICA, 2004), Neural Comp (2006)

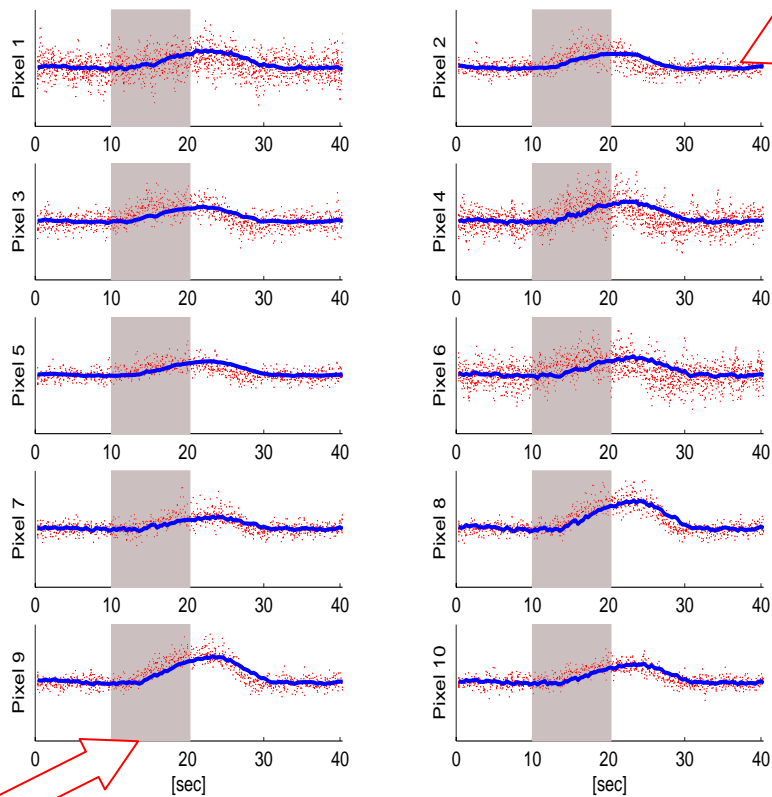
# Instantaneous mixing: Components 1,2 (30%+30% of variance)



# Instantaneous mixing cross correlation between components



# Convolutive mixing (60% of total variance)



Conclusion:

Convolutive ICA is the appropriate mixing model for fMRI

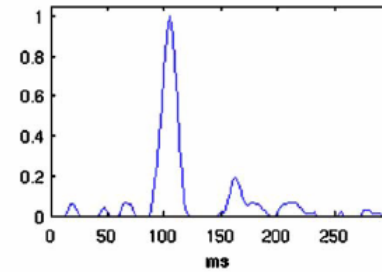
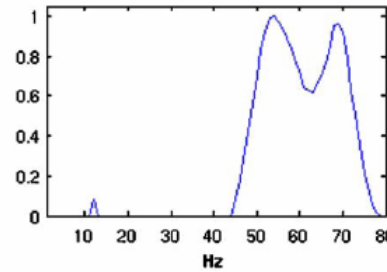
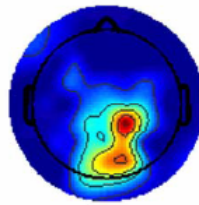
# Data represented as multiway arrays

$$\begin{array}{c}
 \text{[Matrix]} = \sum_{\lambda=1}^F \mathbf{a}_{\lambda} \mathbf{s}_{\lambda} \\
 x_{i_1 i_2} = \sum_{\lambda=1}^F a_{i_1 \lambda} s_{i_2 \lambda} + e_{i_1 i_2} \\
 \text{Factor Analysis}
 \end{array}
 \quad
 \begin{array}{c}
 \text{[3-way array]} = \sum_{\lambda=1}^F \mathbf{a}_{\lambda} \mathbf{d}_{\lambda} \mathbf{s}_{\lambda} \\
 x_{i_1 i_2 i_3} = \sum_{\lambda=1}^F a_{i_1 \lambda} d_{i_2 \lambda} s_{i_3 \lambda} + e_{i_1 i_2 i_3} \\
 \text{PARAFAC}
 \end{array}$$

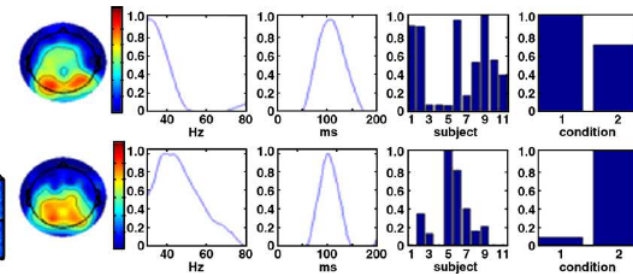
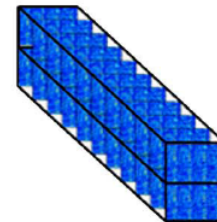
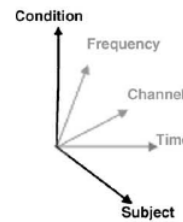
EEG visual response to meaningful vs non-meaningful drawings (N=11).

Fig. 1. Graphical representation of the factor analysis to the left and the PARAFAC decomposition of a 3-way array to the right. Like the factor analysis, PARAFAC decomposes the data into factor effects pertaining to each modality.  $F$  denotes the number of factors.

3-way analysis:  
Channel\*freq\*time



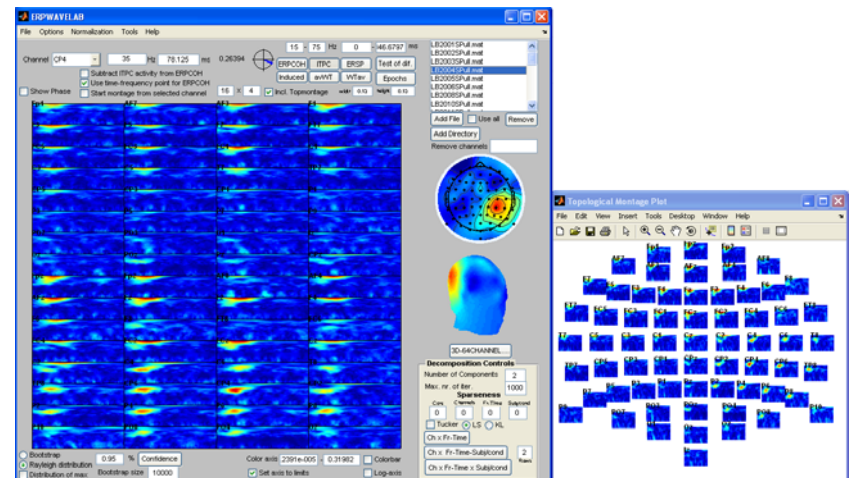
5-way analysis:  
Channel\*freq\*time\*subject\*condition



Mørup et al. NeuroImage (2005), NeuroImage (2008)

# ERPWAVELAB

- Interfaced with EEGLAB
- Single subject analysis
  - Artifact rejection in the time/freq domain
  - NMF decomposition
  - Cross coherence tracking
- Multi subject analysis
  - Clustering
  - Analysis of Variance (ANOVA)
  - Tensor decomposition



Toolbox download from [www.erpwavelab.com](http://www.erpwavelab.com)

Mørup et al. J. Neuroscience Methods (2007),



# Conclusions

- Explorative methods can supplement hypothesis testing.
- Multivariate models can detect brain-wide activation networks.
- Generalizability/learning curves show surprising features: Retarded learning in large systems, crossing learning curves
- Linear hidden variable models form a flexible family including PCA, FA, ICA, Clustering
- New efficient tools for estimation and optimization of multi-linear models are emerging, applications in group analysis, repeat trial models ( NeuroImage, in press)

# Acknowledgments

Lundbeck Foundation ([www.cimbi.org](http://www.cimbi.org))  
NIH Human Brain Project grant ( P20 MH57180)  
MAPAWAMO / EU Commission  
Danish Research Councils

[www.imm.dtu.dk/cisp](http://www.imm.dtu.dk/cisp)  
[hendrix.imm.dtu.dk](http://hendrix.imm.dtu.dk)

