# Introduction to Independent Component Analysis

Jarmo Hurri, Patrik Hoyer, Aapo Hyvärinen

# Course timetable (this week)

- today: ICA examples and mathematical background

- tomorrow: ICA model, decorrelation, non-Gaussianity, FastICA (me)

- Wednesday: practical considerations, ICA by maximum likelihood, image representations (Patrik)

# Course timetable (next week)

- next Monday: nonnegative sparse coding, modeling residual dependencies (Patrik)

- next Tuesday: recent advances and open questions (Aapo)

- <u>all classes</u> at 14:15–16:00 in this room (A414)
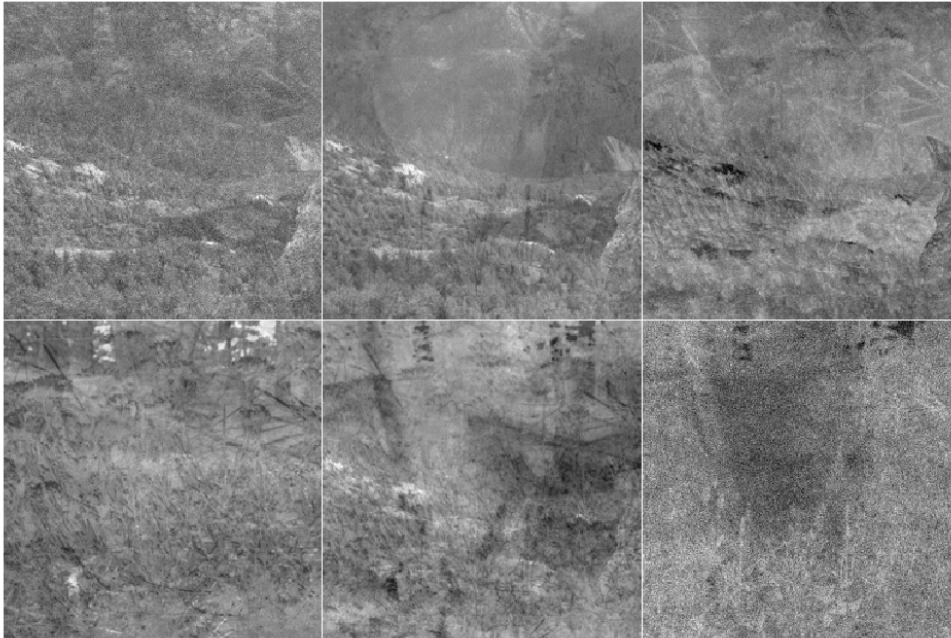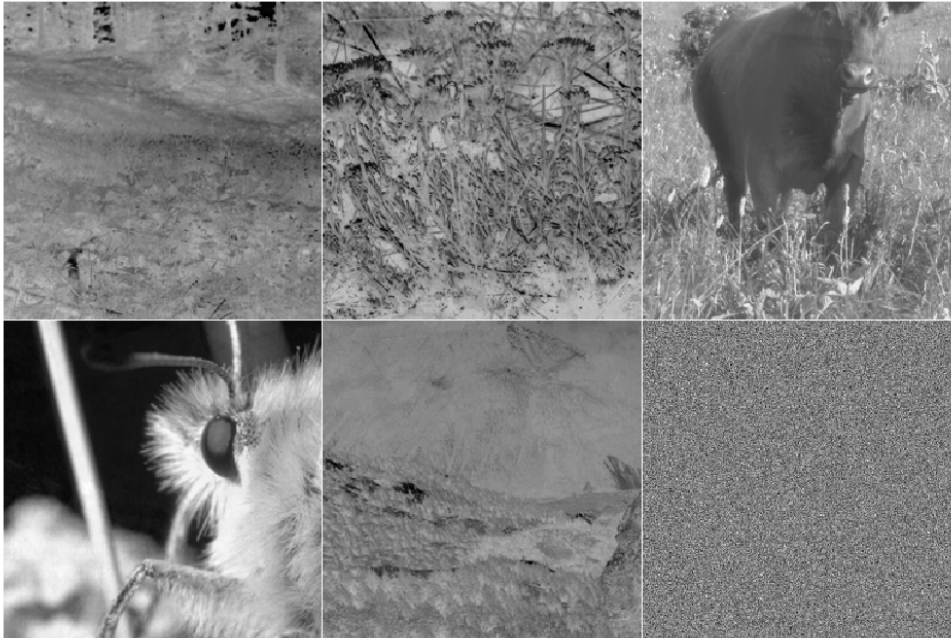
- `http://www.cs.helsinki.fi/jarmo.hurri`

# Contents

# Example



- 6 images

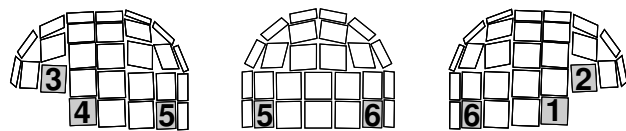- linear mixtures of 6 originals

- determine originals

# Independent components

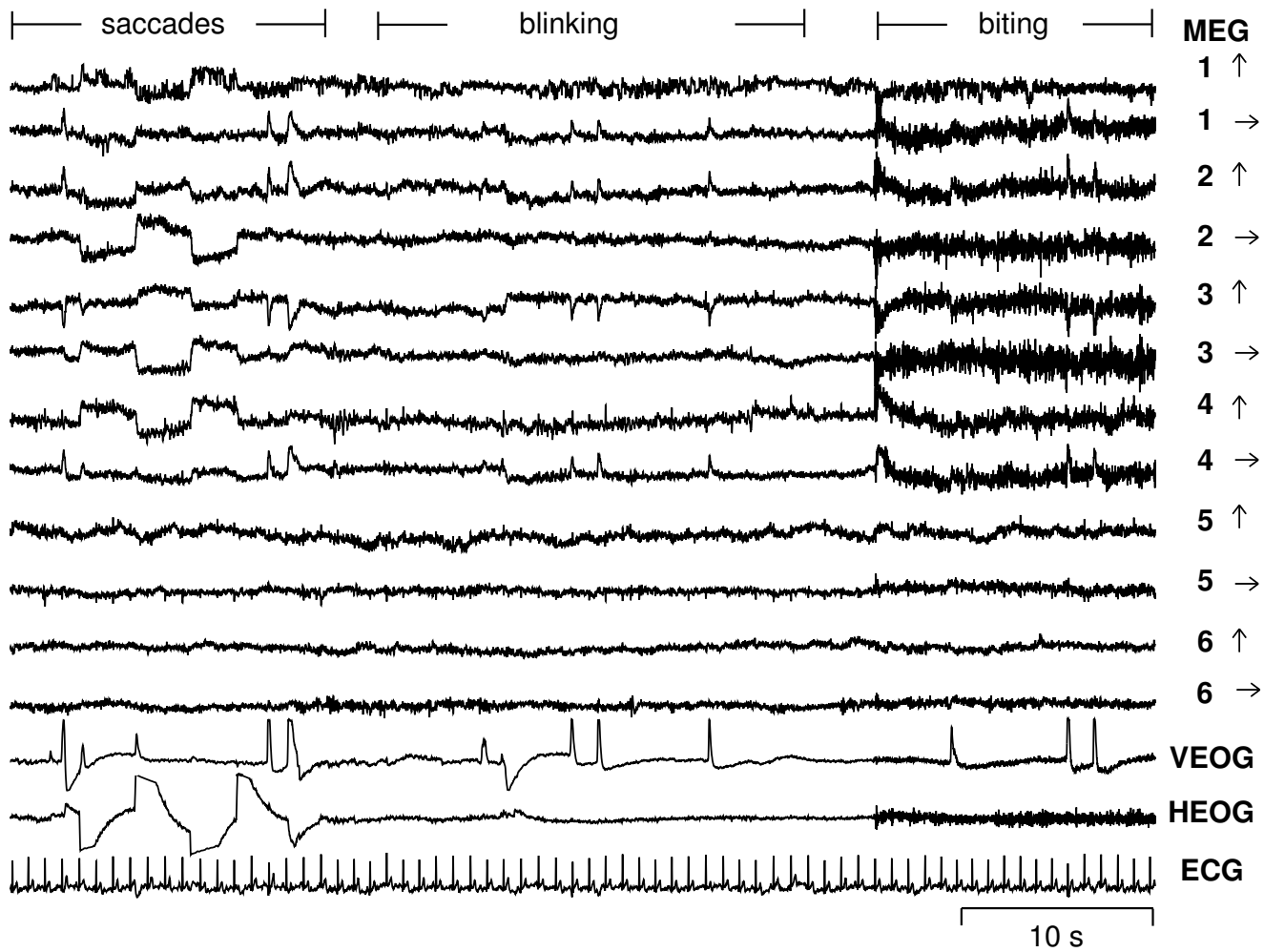# Independent components



- independent latent (hidden) variables

- linear phenomenon

MEG ⎡ 1000 fT/cm

EOG ⎡ 500 µV

ECG ⎡ 500 µV

├— saccades —┤  ├— blinking —┤  ├— biting —┤   **MEG**

**1** ↑

**1** →

**2** ↑

**2** →

**3** ↑

**3** →

**4** ↑

**4** →

**5** ↑

**5** →

**6** ↑

**6** →

**VEOG**

**HEOG**

**ECG**

├— 10 s —┤

IC1

IC2

IC3

IC4

IC5

IC6

IC7

IC8

IC9

10 s

# Some ICA application areas

- biomedical signal analysis (EEG/MEG/MRI/fMRI)

- computational neuroscience

- multispectral image analysis

- bioinformatics (transcriptome analysis)

- gas sensor array analysis

- telecommunications

# II.Multivariate optimization (1/2)

- differentiable objective $f(\mathbf{w})$

- derivative: local linear approximation of change $\Delta f$

- $f'(\mathbf{w}) = [\partial f / \partial w_1 \ \partial f / \partial w_2 \ \cdots \ \partial f / \partial w_n]$

- $f(\mathbf{w} + \Delta \mathbf{w}) - f(\mathbf{w}) \approx f'(\mathbf{w}) \Delta \mathbf{w}$

- gradient rule: $\mathbf{w}(k + 1) = \mathbf{w}(k) \pm \alpha f'(\mathbf{w})^T$

# Multivariate optimization (2/2)

- differentiable constraint $h(\mathbf{w}) = 0$

- necessary conditions:

$$f'(\mathbf{w}) + \lambda h'(\mathbf{w}) = \mathbf{0}^T$$

- different iterative methods (e.g., projected gradient)

# III. Multivariate statistics

- notation

- basic statistics

- multivariate Gaussian density

- principal component analysis

- statistical independence

# Notation

- random variables: $x, y, \ldots$

- density functions: $p_x(x), p_y(y), \ldots$

- random vectors: $\mathbf{x} = [x_1 \ x_2 \ \cdots \ x_n]^T$

  - each component a (continuous) random variable
  - $F(\mathbf{x}_0) = \mathsf{P}(\mathbf{x} \le \mathbf{x}_0)$
  - $p_{\mathbf{x}}(\mathbf{x}_0) = \frac{\partial}{\partial x_1} \frac{\partial}{\partial x_2} \cdots \frac{\partial}{\partial x_n} F(\mathbf{x}) \Big|_{\mathbf{x}=\mathbf{x}_0}$
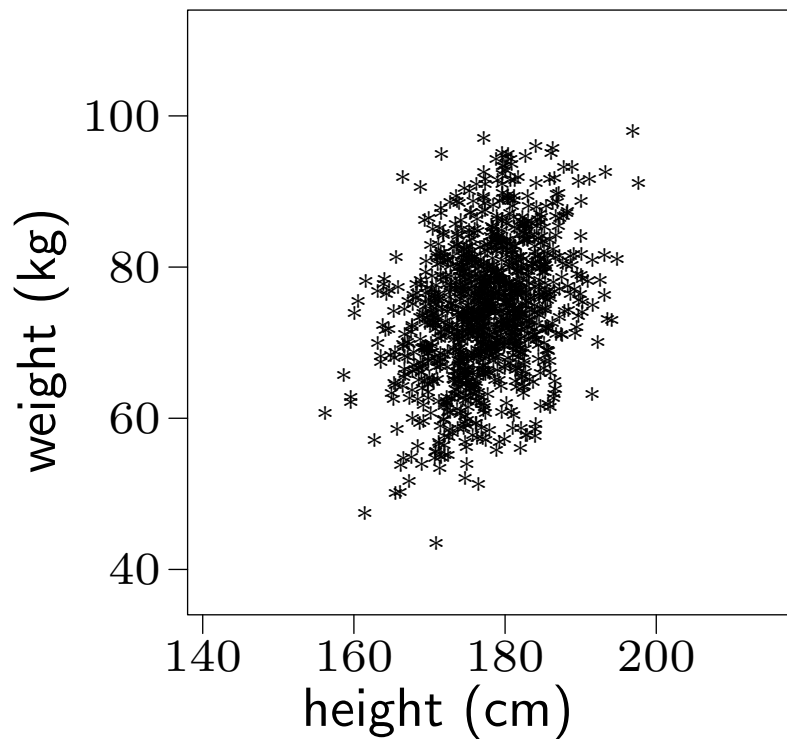
# Basic statistics (1/3)

- expectation: $\mathsf{E}\{g(\mathbf{x})\} = \int g(\mathbf{x})p_{\mathbf{x}}(\mathbf{x})\,d\mathbf{x}$

- mean: $\mathsf{E}\{\mathbf{x}\} = [\mathsf{E}\{x_1\}\ \ \mathsf{E}\{x_2\}\ \ \cdots\ \ \mathsf{E}\{x_n\}]^T$

- correlation matrix:

$$\mathbf{R}_{\mathbf{x}} = \mathsf{E}\{\mathbf{x}\mathbf{x}^T\} = \begin{bmatrix} \mathsf{E}\{x_1^2\} & \mathsf{E}\{x_1 x_2\} & \cdots & \mathsf{E}\{x_1 x_n\} \\ \mathsf{E}\{x_2 x_1\} & \mathsf{E}\{x_2^2\} & \cdots & \mathsf{E}\{x_2 x_n\} \\ \vdots & \vdots & \ddots & \vdots \\ \mathsf{E}\{x_n x_1\} & \mathsf{E}\{x_n x_2\} & \cdots & \mathsf{E}\{x_n^2\} \end{bmatrix}$$

# Basic statistics (2/3)

- covariance matrix $\mathbf{C_x} = \mathsf{E}\left\{ (\mathbf{x} - \mathsf{E}\{\mathbf{x}\})(\mathbf{x} - \mathsf{E}\{\mathbf{x}\})^T \right\}$

  - note: if $\mathsf{E}\{\mathbf{x}\} = \mathbf{0}$, then $\mathbf{C_x} = \mathbf{R_x}$
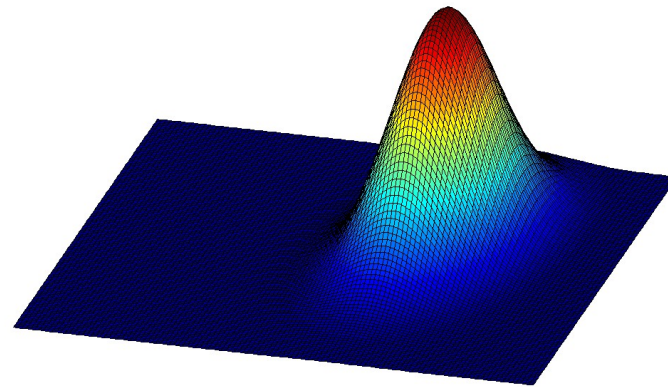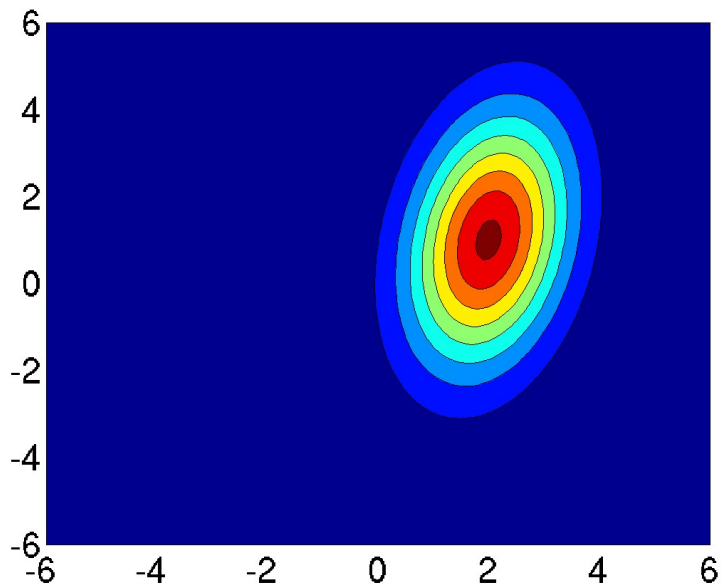
# Basic statistics (3/3)



- $x_1$: height, $x_2$: weight

- $\mathbf{C_x} = \begin{bmatrix} 37.26 & 17.95 \\ 17.95 & 77.32 \end{bmatrix}$

- $\text{cov}\{x_1, x_2\} > 0$

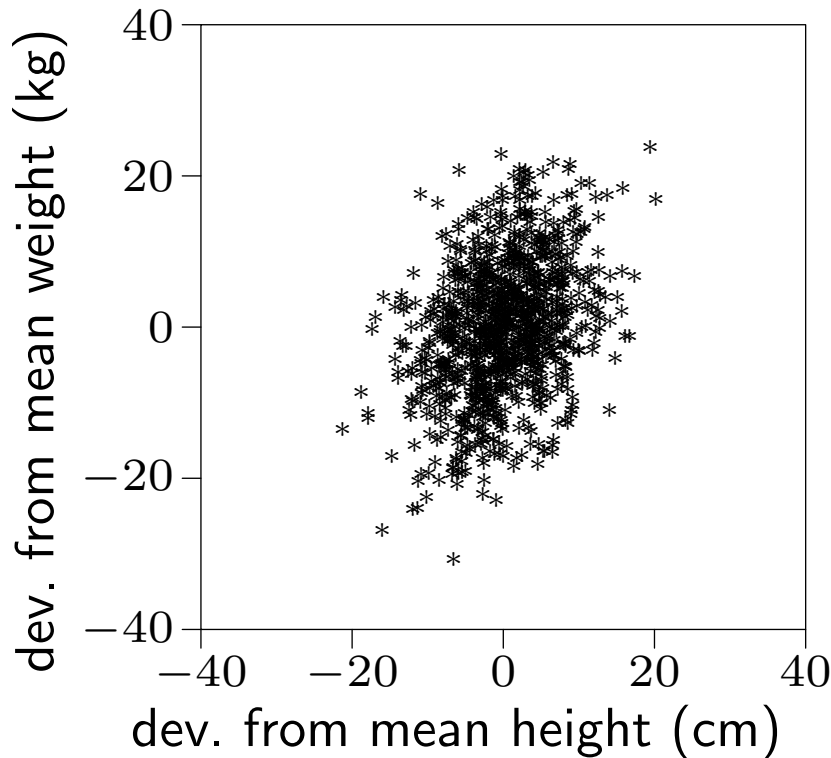- $\text{var}\{x_2\} > \text{var}\{x_1\}$

# Multivariate Gaussian density (1/2)

# Multivariate Gaussian density (2/2)

- $p_{\mathbf{x}}(\mathbf{x}) = K \exp\left(-\frac{1}{2}(\mathbf{x} - \mathbf{m_x})^T \mathbf{C_x}^{-1}(\mathbf{x} - \mathbf{m_x})\right)$

- $K = \left((2\pi)^{n/2} \det(\mathbf{C_x})^{1/2}\right)^{-1}$

- $\mathsf{E}\{\mathbf{x}\} = \mathbf{m_x}$

- covariance matrix $\mathbf{C_x}$

- completely specified by 1st and 2nd order statistics

# Principal component analysis (1/4)



- describe data with <u>one</u> linear projection

- $\mathbf{w} = \begin{bmatrix} w_1 \ w_2 \ \cdots \ w_n \end{bmatrix}^T$, $\|\mathbf{w}\|^2 = 1$

- new data $\mathbf{x}_* = \left( \mathbf{w}^T \mathbf{x} \right) \mathbf{w}$

- min $f(\mathbf{w}) = \mathsf{E} \left\{ \|\mathbf{x} - \mathbf{x}_*\|^2 \right\}$, s.t. $h(\mathbf{w}) = \|\mathbf{w}\|^2 - 1 = 0$

# Principal component analysis (2/4)

$$\mathbf{u} \colon \mathbf{u}^T\mathbf{w} = 0, \|\mathbf{u}\|^2 = 1$$

$$
\begin{aligned}
f(\mathbf{w}) &= \mathsf{E}\{\|\underbrace{(\mathbf{w}^T\mathbf{x})\mathbf{w} + (\mathbf{u}^T\mathbf{x})\mathbf{u}}_{=\mathbf{x}} - \underbrace{(\mathbf{w}^T\mathbf{x})\mathbf{w}}_{=\mathbf{x}_*}\|^2\} \\
&= \mathsf{E}\left\{\|(\mathbf{u}^T\mathbf{x})\mathbf{u}\|^2\right\} \text{ expected norm of proj.} \\
&= \mathsf{E}\left\{\mathbf{u}^T(\mathbf{u}^T\mathbf{x})(\mathbf{u}^T\mathbf{x})\mathbf{u}\right\} \\
&= \mathbf{u}^T\mathsf{E}\left\{\mathbf{x}\mathbf{x}^T\right\}\mathbf{u} = \mathbf{u}^T\mathbf{R_x}\mathbf{u}
\end{aligned}
$$

# Principal component analysis (3/4)

$$f(\mathbf{u}) = \mathbf{u}^T \mathbf{R_x} \mathbf{u}$$
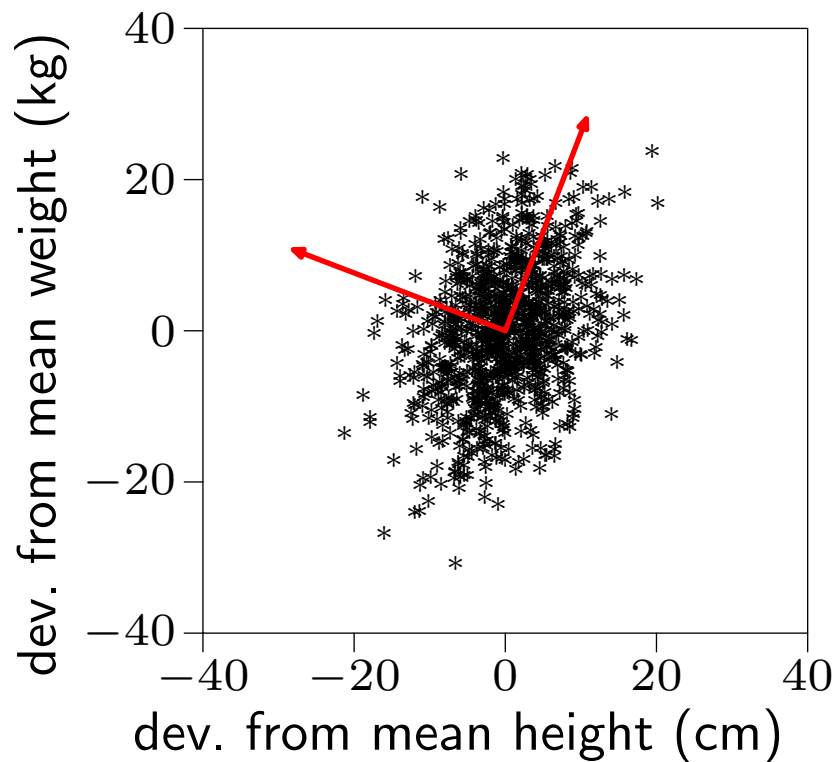
$$h(\mathbf{u}) = \mathbf{u}^T \mathbf{u} - 1 = 0$$

$$f'(\mathbf{u}_{\mathsf{opt}}) - \lambda h'(\mathbf{u}_{\mathsf{opt}}) = \mathbf{0}^T$$

$$2\mathbf{u}_{\mathsf{opt}}^T \mathbf{R_x} - \lambda 2\mathbf{u}_{\mathsf{opt}}^T = \mathbf{0}^T$$

$$\mathbf{R_x} \mathbf{u}_{\mathsf{opt}} = \lambda \mathbf{u}_{\mathsf{opt}}$$

$$f(\mathbf{u}_{\mathsf{opt}}) = \lambda$$

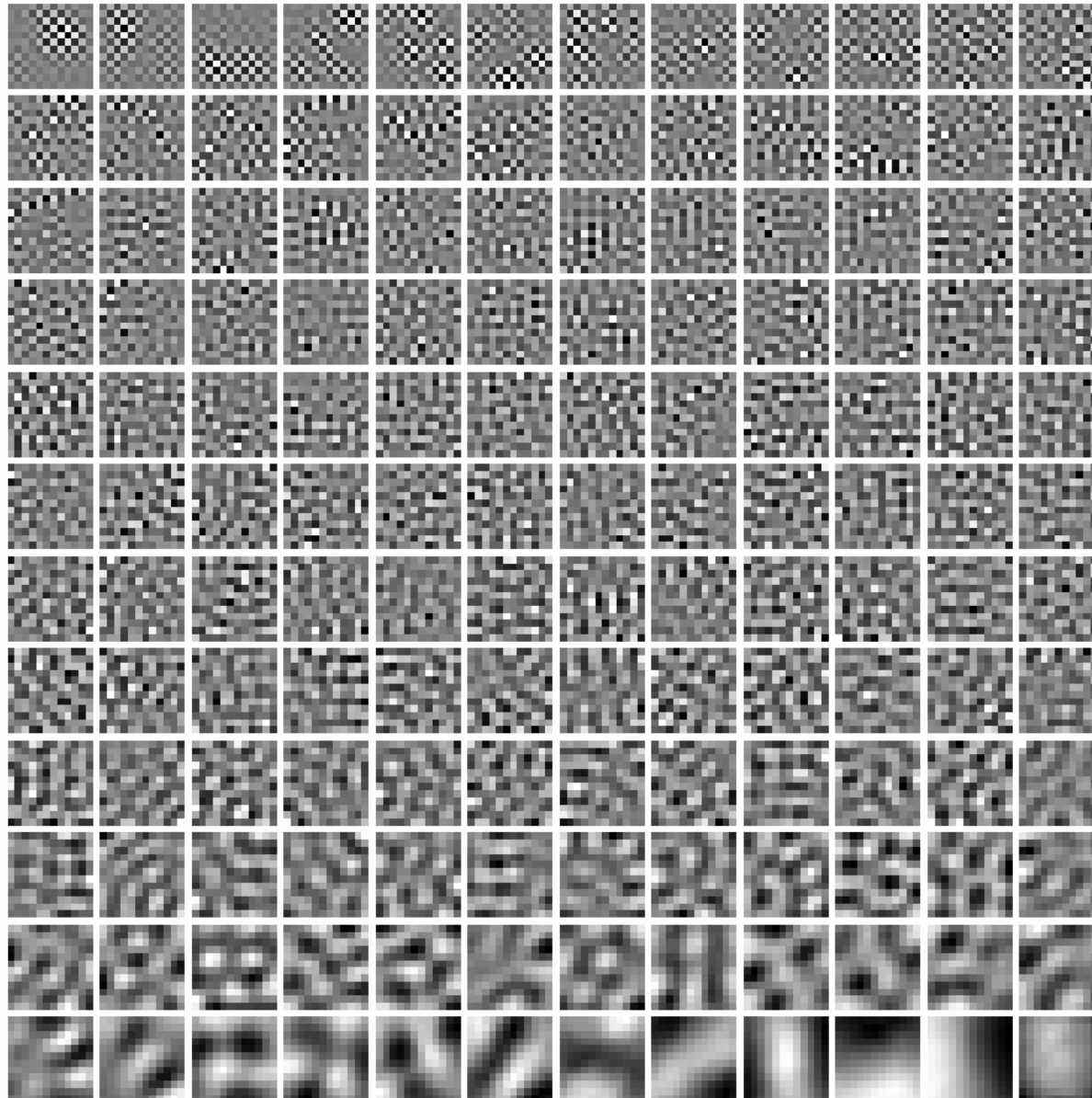# Principal component analysis (4/4)



- orthogonal PCA basis

- deflationary     minima     /
  maxima of variance

- dimensionality reduction

- noise attenuation

# PCA example — original

# PCA example — PCA basis

# PCA example — 90% compressed

# PCA example — 75% compressed

# PCA example — 50% compressed

# PCA example — original

# Statistical independence

- random variables $x$ and $y$

- statistical independence: knowing the value of $x$ does not provide information about the distribution of $y$

- $p_y\left(y|x\right) = \frac{p_{x,y}(x,y)}{p_x(x)} = p_y\left(y\right)$
  $\Leftrightarrow p_{x,y}\left(x,y\right) = p_x\left(x\right)p_y\left(y\right)$

# IV. Estimation theory

- basic concepts

- maximum-likelihood estimation

# Basic concepts (1/3)

- <u>estimation</u>: finding an approximate value or distribution of some <u>parameter</u> (e.g., mean, standard deviation) from random samples of the population

- <u>estimator</u>: a function computing the approximate value or distribution from the samples (e.g., $\hat{\mu} = 1/T \sum_{i=1}^{T} x(i)$)

- <u>estimate</u>: numerical value of the estimator for a given set of samples

# Basic concepts (2/3)

- some good properties of estimators:

  - unbiasedness: $\mathsf{E}\left\{\hat{\theta}\right\} = \theta$

  - consistency:

$$\forall \epsilon > 0 \left[ \lim_{T \to \infty} \mathsf{P}\left( \left\| \theta - \hat{\theta} \right\| < \epsilon \right) = 1 \right]$$

# Basic concepts (3/3)

- **data model**: mathematical representation of data

  - parametric / nonparametric
  - static / dynamic
  - probabilistic / deterministic
  - example: $p_x\left(x|\mu,\sigma\right) = K e^{-\frac{(x-\mu)^2}{2\sigma^2}}$

# Maximum-likelihood estimation (1/3)

- select data model parameter values that maximize the probability of the observed data

- $\hat{\theta}_{\mathsf{ML}} = \arg\max_\theta \mathsf{P}\left(x(1), x(2), ..., x(T) \,|\, \theta\right)$

- continuous data: solve the likelihood equation

$$\left.\frac{\partial}{\partial \theta} \ln p_{x(1),...,x(T)}\left(x(1), ..., x(T)|\theta\right)\right|_{\theta=\hat{\theta}_{\mathsf{ML}}} = 0$$

# Maximum-likelihood estimation (2/3)

- when samples are independent

$$p_{x(1),...,x(T)}\left(x(1),...,x(T)|\theta\right) = \prod_{i=1}^{T} p_{x(i)}\left(x(i)|\theta\right)$$

- example: mean of a Gaussian

- $p_{x(1),...,x(T)}\left(x(1),...,x(T)|\mu\right) = K^{T} \prod_{i=1}^{T} e^{-\frac{(x(i)-\mu)^2}{2\sigma^2}}$

# Maximum-likelihood estimation (3/3)

$$\frac{\partial}{\partial \mu} \left[ -\frac{1}{2\sigma^2} \sum_{i=1}^{T} \left( x(i) - \mu \right)^2 \right] \Bigg|_{\mu = \mu_{\mathsf{ML}}} = 0$$

$$\sum_{i=1}^{T} \left( x(i) - \mu_{\mathsf{ML}} \right) = 0$$

$$\mu_{\mathsf{ML}} = \frac{1}{T} \sum_{i=1}^{T} x(i)$$

# Tomorrow

ICA in images and video (and equations...)