

Title: Informal Discussion of Eternal Inflation and String Theory

Date: Feb 01, 2005 02:00 PM

URL: <http://pirsa.org/05020001>

Abstract:

measures in all sciences:

measure 2 quantities A, B

⇒ are they related?

⇒ how can one quantify degree?

temperature at town X } ✓
rainfall -u- }

temperature at Waterloo } ?
chinese stock market index }

births } ?
of breeding storks }

consumption of drug } ??
health stats }

Applications:

clustering of objects

decomposition of signals into
~ independent components

classification

⋮

clustering of objects

decomposition of signals into
~ independent components

classification

⋮

clustering of objects

decomposition of signals into
~ independent components

classification

⋮

clustering of objects

decomposition of signals into
~ independent components

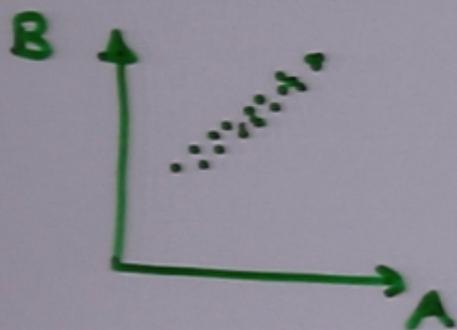
classification

⋮

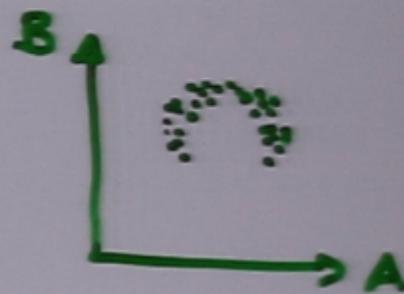
Correlation coefficient:

$$C_{AB} = \langle AB \rangle - \langle A \rangle \langle B \rangle$$

only linear relationship!



o.k.



??

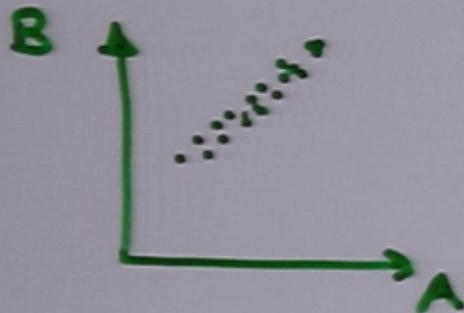
Universal,
information - theoretic
measure of dependencies:

Mutual Information

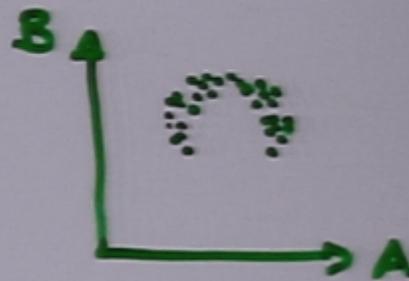
Correlation coefficient :

$$C_{AB} = \langle AB \rangle - \langle A \rangle \langle B \rangle$$

only linear relationship!



o.k.



??

Universal,
information - theoretic
measure of dependencies:

Mutual Information



o.k.



??

Universal,
information - theoretic
measure of dependencies:

Mutual Information

Shannon Theory

algorithmic (Kolmogorov)
inform. Theory

Correlation coefficient =

$$C_{AB} = \langle AB \rangle - \langle A \rangle \langle B \rangle$$

only linear relationships!



O.K.



??

Shannon information theory:

$X, Y =$ random variables

If X, Y discrete:

$$H(X) = -\sum_i p(x_i) \log p(x_i)$$

Shannon entropy

$$I(X, Y) = H(X) + H(Y) - H(X, Y)$$

M.I.

$$= H(X) - H(X|Y)$$

$$= \sum_{ij} p_{ij} \log \frac{p_{ij}}{p_i p_j}$$

If X, Y continuous:

$X, Y =$ random variables

If X, Y discrete:

$$H(X) = -\sum_i p(x_i) \log p(x_i)$$

Shannon entropy

$$I(X, Y) = H(X) + H(Y) - H(X, Y)$$

M.I.

$$= H(X) - H(X|Y)$$

$$= \sum_{i,j} p_{ij} \log \frac{p_{ij}}{p_i p_j}$$

If X, Y continuous:

- make binning (coarse graining)

$$H(X) = - \sum_i p(x_i) \log p(x_i)$$

Shannon entropy

$$I(X, Y) = H(X) + H(Y) - H(X, Y)$$

M.I.

$$= H(X) - H(X|Y)$$

$$= \sum_{i,j} p_{ij} \log \frac{p_{ij}}{p_i p_j}$$

If X, Y continuous:

- make binning (coarse graining)
- take bin size $\rightarrow 0$:

$$I(X, Y) = \int dx dy p(x, y) \log \frac{p(x, y)}{p(x) p(y)}$$

$$I(X, Y) = H(X) + H(Y) - H(X, Y)$$

M.C.

$$= H(X) - H(X|Y)$$

$$= \sum_{i,j} p_{ij} \log \frac{p_{ij}}{p_i p_j}$$

If X, Y continuous:

- make binning (coarse graining)
- take bin size $\rightarrow 0$:

$$I(X, Y) = \int dx dy p(x, y) \log \frac{p(x, y)}{p(x) p(y)}$$

If X, Y discrete:

$$H(X) = -\sum_i p(x_i) \log p(x_i)$$

Shannon entropy

$$I(X, Y) = H(X) + H(Y) - H(X, Y)$$

M.I.

$$= H(X) - H(X|Y)$$

$$= \sum_{i,j} p_{ij} \log \frac{p_{ij}}{p_i p_j}$$

If X, Y continuous:

- make binning (coarse graining)
- take bin size $\rightarrow 0$:

Shannon entropy

$$I(X, Y) = H(X) + H(Y) - H(X, Y)$$

M.E.

$$= H(X) - H(X|Y)$$

$$= \sum_{i,j} p_{ij} \log \frac{p_{ij}}{p_i p_j}$$

If X, Y continuous:

- make binning (coarse graining)
- take bin size $\rightarrow 0$:

$$I(X, Y) = \int dx dy p(x, y) \log \frac{p(x, y)}{p(x) p(y)}$$

$X, Y =$ random variables

If X, Y discrete:

$$H(X) = -\sum_i p(x_i) \log p(x_i)$$

Shannon entropy

$$I(X, Y) = H(X) + H(Y) - H(X, Y)$$

M.I.

$$= H(X) - H(X|Y)$$

$$= \sum_{i,j} p_{ij} \log \frac{p_{ij}}{p_i p_j}$$

If X, Y continuous:

• make binning (coarse graining)

If X, Y discrete:

$$H(X) = -\sum_i p(x_i) \log p(x_i)$$

Shannon entropy

$$I(X, Y) = H(X) + H(Y) - H(X, Y)$$

M.I.

$$= H(X) - H(X|Y)$$

$$= \sum_{i,j} p_{ij} \log \frac{p_{ij}}{p_i p_j}$$

If X, Y continuous:

- make binning (coarse graining)
- take bin size $\rightarrow 0$:

Shannon entropy

$$I(X, Y) = H(X) + H(Y) - H(X, Y)$$

M.T.

$$= H(X) - H(X|Y)$$

$$= \sum_{i,j} p_{ij} \log \frac{p_{ij}}{p_i p_j}$$

If X, Y continuous:

- make binning (coarse graining)
- take bin size $\rightarrow 0$:

$$I(X, Y) = \int dx dy p(x, y) \log \frac{p(x, y)}{p(x) p(y)}$$

N.B. : $H(x) \rightarrow \infty$ when bin size $\rightarrow 0$!

Properties of MI :

1.) $I(x,y) = I(y,x)$ symm.

2.) $I(x,y) \geq 0$

$= 0$ only if X, Y indep.
 $p(x,y) = p(x) \cdot p(y)$

3.) $1 - \frac{I(x,y)}{H(x,y)}$ is distance
(Δ -inequality)

4.) $I(x,y)$ invar under
homeomorphisms

$X \rightarrow \phi(x)$ ϕ^{-1}, ψ^{-1} exist

$Y \rightarrow \psi(y)$

$H(x), H(x,y) !$

$X, Y =$ random variables

If X, Y discrete:

$$H(X) = -\sum_i p(x_i) \log p(x_i)$$

Shannon entropy

$$I(X, Y) = H(X) + H(Y) - H(X, Y)$$

M.I.

$$= H(X) - H(X|Y)$$

$$= \sum_{i,j} p_{ij} \log \frac{p_{ij}}{p_i p_j}$$

If X, Y continuous:

If X, Y discrete:

$$H(X) = - \sum_i p(x_i) \log p(x_i)$$

Shannon entropy

$$I(X, Y) = H(X) + H(Y) - H(X, Y)$$

M.I.

$$= H(X) - H(X|Y)$$

$$= \sum_{i,j} p_{ij} \log \frac{p_{ij}}{p_i p_j}$$

NB. : $H(x) \rightarrow \infty$ when bin size $\rightarrow 0$

(12
6
!

Properties of MI:

- 1) $I(x, y) = I(y, x)$ symm.
- 2) $I(x, y) \geq 0$
 $= 0$ only if X, Y indep.
 $p(x, y) = p(x) \cdot p(y)$
- 3) $1 - \frac{I(x, y)}{H(x, y)}$ is distance
(Δ -inequality)
- 4) $I(x, y)$ invar under homeomorphisms

$X \rightarrow \phi(x)$ ϕ^{-1}, ψ^{-1} exist

2.) $I(x,y) \geq 0$

$= 0$

only if X, Y indep.

$p(x,y) = p(x) \cdot p(y)$

3) $1 - \frac{I(x,y)}{H(x,y)}$

is distance

(Δ -inequality)

4) $I(x,y)$ invar under
homeomorphisms

$X \rightarrow \phi(x)$

$Y \rightarrow \psi(y)$

ϕ^{-1}, ψ^{-1} exist

(not true for $H(x), H(x,y)$!)

5) for fixed covar matrix C :

$I(x,y)$ is minimal for Gaussians

$$2.) I(x, Y) \geq 0$$

$$= 0$$

only if X, Y indep.

$$p(x, y) = p(x) \cdot p(y)$$

$$3.) 1 - \frac{I(x, Y)}{H(x, Y)} \text{ is distance}$$

(Δ -inequality)

4) $I(x, Y)$ invar under homeomorphisms

$$X \rightarrow \phi(x) \quad \phi^{-1}, \psi^{-1} \text{ exist}$$

$$Y \rightarrow \psi(y)$$

(not true for $H(x), H(x, Y)$!)

5) for fixed covar matrix C :

$I(x, Y)$ is minimal for Gaussian

$$2.) I(x, Y) \geq 0$$

$$= 0$$

only if X, Y indep.

$$p(x, y) = p(x) \cdot p(y)$$

$$3.) 1 - \frac{I(x, Y)}{H(x, Y)} \text{ is distance}$$

(Δ -inequality)

4) $I(x, Y)$ invar under
homeomorphisms

$$X \rightarrow \phi(x)$$

$$Y \rightarrow \psi(y)$$

ϕ^{-1}, ψ^{-1} exist

(not true for $H(x), H(x, Y)$!)

5) for fixed covar matrix C :

$I(x, Y)$ is minimal for Gaussian

$$\Rightarrow I(X_1, X_2, \dots, X_m) \geq \frac{1}{2} \log \frac{\det C}{C_{11} C_{22} \dots C_{mm}}$$

MI for ≥ 3 Variables:

- there are various definitions in the literature !

ICA literature :

$$I(x_1 \dots x_n) = H(x_1) + \dots + H(x_n) - H(x_1 \dots x_n)$$

Grouping property :

def.: $W = (X, Y)$: joint random variable

$$I(x, y, z, \dots) = I(W, z, \dots) + I(x, y)$$

MI for ≥ 3 Variables:

- there are various definitions in the literature !

ICA literature :

$$I(x_1 \dots x_n) = H(x_1) + \dots + H(x_n) - H(x_1 \dots x_n)$$

Grouping property :

def.: $W = (X, Y)$: joint random variable

$$I(x, y, z, \dots) = I(W, z, \dots) + I(x, y)$$

$$I(X_1 \dots X_n) = H(X_1) + \dots + H(X_n) \\ - H(X_1 \dots X_n)$$

Grouping property :

def.: $W = (X, Y)$: joint random variable

$$I(X, Y, Z, \dots) = I(W, Z, \dots) \\ + I(X, Y)$$

\Rightarrow hierarchical cluster decomposition of MI:

$$MI = MI \text{ within clusters} \\ + \\ MI \text{ between clusters}$$

$$I(X_1, \dots, X_n) = H(X_1) + \dots + H(X_n) - H(X_1, \dots, X_n)$$

Grouping property :

def.: $W = (X, Y)$: joint random variable

$$I(X, Y, Z, \dots) = I(W, Z, \dots) + I(X, Y)$$

\Rightarrow hierarchical cluster decomposition of MI:

$$MI = \underset{\substack{\text{MI within clusters} \\ + \\ \text{MI between clusters}}}{MI}$$

$X = x_1 x_2 x_3 \dots x_N$: symbol string
 $x_i \in$ "alphabet" \mathcal{A}

U = universal (Turing-) computer,
i.e. PC

$C_u(X)$ = length (in bits) of shortest
program which prints X on U
and then halts U

"Kolmogorov complexity of X "

Thus:

- $\lim_{N \rightarrow \infty} \frac{1}{N} C_u(X)$ is indep. of U
for nearly all X

- if X is randomly drawn from
stochastic process with entropy H
per letter:

Algorithmic Inform. Theory:

$X = x_1 x_2 x_3 \dots x_N$: symbol string
 $x_i \in$ "alphabet" \mathcal{A}

U = universal (Turing-) computer,
i.e. PC

$C_u(X)$ = length (in bits) of shortest
program which prints X on U
and then halts U

"Kolmogorov complexity of X "

Thus:

- $\lim_{N \rightarrow \infty} \frac{1}{N} C_u(X)$ is indep. of U
for nearly all X

- if X is randomly drawn from

U = universal (Turing-) computer,
i.e. PC

$C_u(X)$ = length (in bits) of shortest
program which prints X on U
and then halts U

"Kolmogorov complexity of X "

Thus:

- $\lim_{N \rightarrow \infty} \frac{1}{N} C_u(X)$ is indep. of U
for nearly all X

- if X is randomly drawn from
stochastic process with entropy H
per letter:

$$\Rightarrow \underline{\lim_{N \rightarrow \infty} \frac{1}{N} C_u(X) = H}$$

$C_u(X)$ = length (in bits) of shortest program which prints X on U and then halts U

"Kolmogorov complexity of X "

Thus:

- $\lim_{N \rightarrow \infty} \frac{1}{N} C_u(X)$ is indep. of U for nearly all X

- if X is randomly drawn from stochastic process with entropy H per letter:

$$\Rightarrow \underline{\lim_{N \rightarrow \infty} \frac{1}{N} C_u(X) = H}$$

Algorithmic Inform. Theory:

$X = x_1 x_2 x_3 \dots x_N$: symbol string
 $x_i \in$ "alphabet" \mathcal{A}

U = universal (Turing-) computer,
i.e. PC

$C_u(X)$ = length (in bits) of shortest
program which prints X on U
and then halts U

"Kolmogorov complexity of X "

Thus:

- $\lim_{N \rightarrow \infty} \frac{1}{N} C_u(X)$ is indep. of U
for nearly all X
- if X is randomly drawn from
stochastic process with entropy H

U = universal (Turing-) computer,
i.e. PC

$C_u(X)$ = length (in bits) of shortest
program which prints X on U
and then halts U

"Kolmogorov complexity of X "

Thus:

- $\lim_{N \rightarrow \infty} \frac{1}{N} C_u(X)$ is indep. of U
for nearly all X

- if X is randomly drawn from
stochastic process with entropy H
per letter:

$$\Rightarrow \underline{\lim_{N \rightarrow \infty} \frac{1}{N} C_u(X) = H}$$

- $C_u(x)$ is not computable, i.e.
 \nexists algorithm for calculating it

Lower estimates on $C(x)$:

given by any compression algor.,

e.g. zip, gzip, compress, bzip2,
 ...

Mutual information between X

$$X = x_1 x_2 \dots x_N$$

$$Y = y_1 \dots y_M :$$

$$I_c(X, Y) = C(X) + C(Y) - C(\underbrace{XY}_{\uparrow})$$

concatenation $x_1 \dots x_N y_1 \dots y_M$

notice: $XY \neq YX$,

$$C(YX) = C(XY) + O(\ln N)$$

- $C_u(x)$ is not computable, i.e.

∄ algorithm for calculating it

Lower estimates on $C(x)$:

given by any compression algor.,

e.g. zip, gzip, compress, bzip2,
...

Mutual information between X

$$X = x_1 x_2 \dots x_N$$

$$Y = y_1 \dots y_M :$$

$$I_c(X, Y) = C(X) + C(Y) - C(\underbrace{XY}_{\uparrow})$$

concatenation $x_1 \dots x_N y_1 \dots y_M$

notice: $XY \neq YX$,

$$C(YX) = C(XY) + O(\ln N)$$

Algorithm for ...

Lower estimates on $C(x)$:

given by any compression algor.,

e.g. zip, gzip, compress, bzip2,
...

Mutual information between X

$$X = x_1 x_2 \dots x_N$$

$$Y = y_1 \dots y_M :$$

$$I_c(X, Y) = C(X) + C(Y) - C(\underbrace{XY}_{\uparrow})$$

concatenation $x_1 \dots x_N y_1 \dots y_M$

notice : $XY \neq YX$,

$$C(XY) = C(YX) + O(\ln N)$$

⇒ all properties of Shannon-MI for

A algorithm for calculating it

Lower estimates on $C(x)$:

given by any compression algor.,

e.g. zip, gzip, compress, bzip2,
...

Mutual information between X

$$X = x_1 x_2 \dots x_N$$

$$Y = y_1 \dots y_M$$

$$I_c(X, Y) = C(X) + C(Y) - C(\overbrace{XY}^{\uparrow})$$

concatenation $x_1 \dots x_N y_1 \dots y_M$

notice: $XY \neq YX$,

$$C(XY) = C(YX) + O(\ln N)$$

Lower estimates on $C(X)$:

given by any compression algor.,

e.g. zip, gzip, compress, bzip2,
...

Mutual information between X

$$X = x_1 x_2 \dots x_N$$

$$Y = y_1 \dots y_M :$$

$$I_c(X, Y) = C(X) + C(Y) - C(\underbrace{XY}_{\uparrow})$$

concatenation $x_1 \dots x_N y_1 \dots y_M$

notice : $XY \neq YX$,

$$C(XY) = C(YX) + O(\ln N)$$

⇒ all properties of Shannon-MI for
discrete variables hold up to terms
 $\sim \ln N$

Main difference between algorithmic & Shannon versions:

$H(X)$ = property of random variable

$C(X)$ = property of single realization

Let X be produced by stochastic process X :

estimation of $C(X)$ via compression algor.

- 4 - $H(X)$: via estimate of statistics, e.g. block frequencies, from X

⇒ basically same estimate

$H(X)$ = property of random variable

$C(X)$ = property of single realization

Let X be produced by stochastic process X :

estimation of $C(X)$ via compression algor.

- 4 - $H(X)$: via estimate of statistics, e.g.
block frequencies,
from X

\Rightarrow basically same estimate

Phylogenetic trees

X
Y
Z
⋮

} mitochondrial genomes
of 34 mammals
($\sim 16,000$ bases each)

$I_c(X, Y)$ = algorithmic Mi between X, Y
obtained with standard comp
(e.g. bzip2)

$$D(X, Y) = 1 - \frac{I_c(X, Y)}{C(X, Y)} = \text{normalized distance}$$

Phylogenetic trees

X
Y
Z
⋮

} mitochondrial genomes
of 34 mammals
($\sim 16,000$ bases each)

$I_c(X, Y)$ = algorithmic MI between X, Y
obtained with standard compressors
(e.g. bzip2)

$$D(X, Y) = 1 - \frac{I_c(X, Y)}{C(X, Y)} = \text{normalized distance}$$

X
Y
Z
...

}

mitochondrial genomes
of 31 mammals

($\sim 16,000$ bases each)

$I_c(X, Y)$ = algorithmic Mi between X, Y
obtained with standard compressors
(e.g. bzip2)

$$D(X, Y) = 1 - \frac{I_c(X, Y)}{C(X, Y)} = \text{normalized distance}$$

X
Y
Z
...

} mitochondrial genomes
of 34 mammals
(~ 16,000 bases each)

$I_c(X, Y)$ = algorithmic Mi between X, Y
obtained with standard compressors
(e.g. bzip2)

$$D(X, Y) = 1 - \frac{I_c(X, Y)}{C(X, Y)} = \text{normalized distance}$$

X
Y
Z
...

} mitochondrial genomes
of 34 mammals

($\sim 16,000$ bases each)

$I_c(X, Y)$ = algorithmic Mi between X, Y
obtained with standard compressors
(e.g. bzip2)

$$D(X, Y) = 1 - \frac{I_c(X, Y)}{C(X, Y)} = \text{normalized distance}$$

X
Y
Z
...

} mitochondrial genomes
of 34 mammals

($\sim 16,000$ bases each)

$I_c(X, Y)$ = algorithmic Mi between X, Y
obtained with standard compressors
(e.g. bzip2)

$$D(X, Y) = 1 - \frac{I_c(X, Y)}{C(X, Y)} = \text{normalized distance}$$

X
Y
Z
...

}

mitochondrial genomes
of 34 mammals
(~ 16,000 bases each)

$I_c(X, Y)$ = algorithmic Mi between X, Y
obtained with standard compression
(e.g. bzip2)

$$D = 1 - \frac{I_c(X, Y)}{C(X, Y)} = \text{normalized distance}$$

X
Y
Z
...

}

mitochondrial genomes
of 34 mammals

($\sim 16,000$ bases each)

$I_c(X, Y)$ = algorithmic Mi between X, Y
obtained with standard compression
(e.g. bzip2)

$$D(X, Y) = 1 - \frac{I_c(X, Y)}{C(X, Y)} = \text{normalized distance}$$

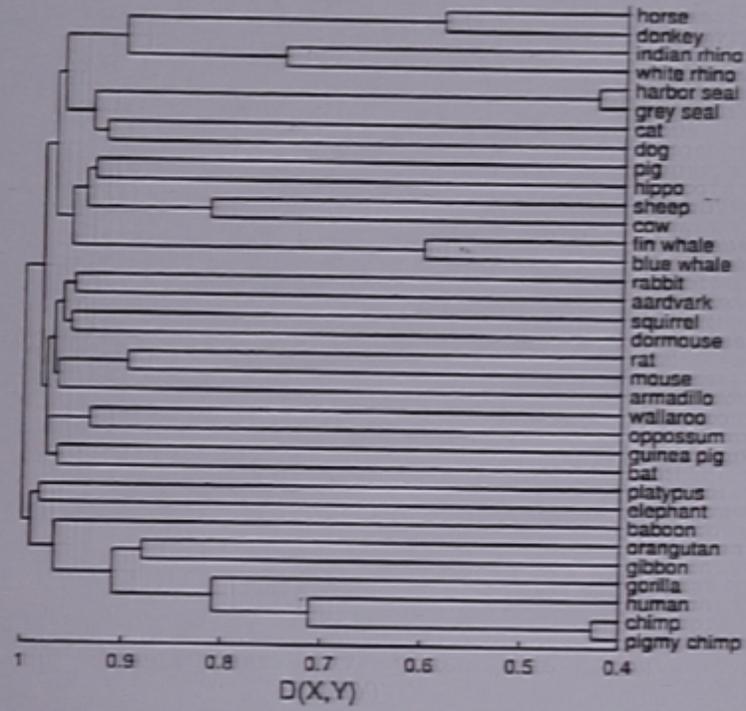


Fig. 1 - Phylogenetic tree for 34 mammals. The heights of nodes are the distances between the joint daughter clusters.

Algorithm :

- ① start with $K = 34$ genomes
- ② form $K \times K$ matrix of distances ←
- ③ search minimal distance, say between species X_j, X_e
- ④ concatenate $X' = X_j X_e$
eliminate species X_j, X_e
&
replace by X'
(thereby $K \rightarrow K - 1$)
- ⑤ if $K > 1$, goto ②

= typical hierarchical clustering alg.,
placement

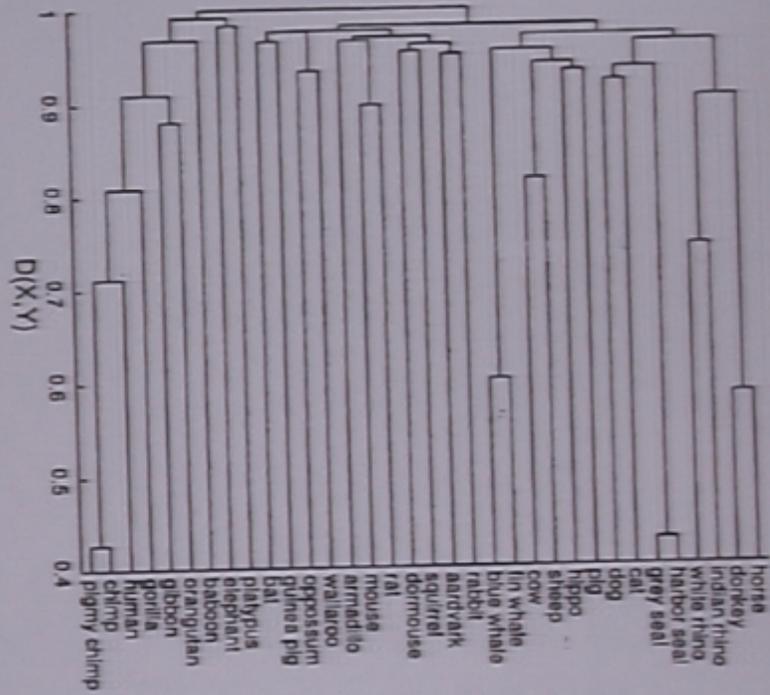


Fig. 1 - Phylogenetic tree for 34 mammals. The heights of nodes are the distances between daughter clusters.

Algorithm:

- ① start with $K = 34$ genomes
- ② form $K \times K$ matrix of distances ←
- ③ search minimal distance, say between species X_j, X_e
- ④ concatenate $X' = X_j X_e$
eliminate species X_j, X_e
&
replace by X'
(thereby $K \rightarrow K-1$)
- ⑤ if $K > 1$, goto ②

= typical hierarchical clustering alg.,
except that replacement

$$X_j, X_e \rightarrow X'$$

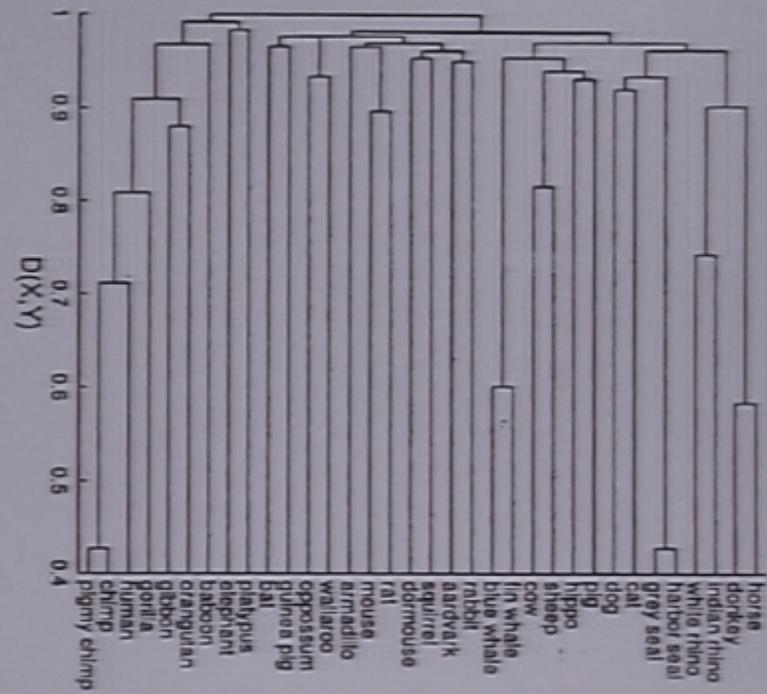


Fig. 1 - Phylogenetic tree for 34 mammals. The heights of nodes are the distances between the joined daughter clusters.

Shannon MI of

Continuous random variables

Estimating ML:

finite sample $(x, y)_1 \dots (x, y)_N$ iid

$\Rightarrow \hat{I}(x, y) \quad ?$

- binning

- fixed grid: very bad!

- adaptive grid

Fraser & Swinney

Darbellay & Vajda

still large systematic error

large statistical error

difficult in high dim's

... but fast!

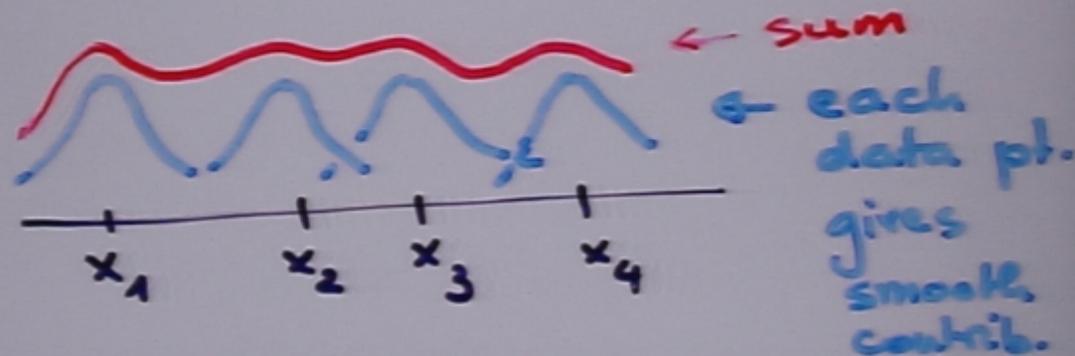
- kernel methods

- adaptive grid

Frazer & Swinney
Darbellay & Vajda

skill! large systematic errors
large statistical --
difficult in high dim's
... but fast!

- kernel methods

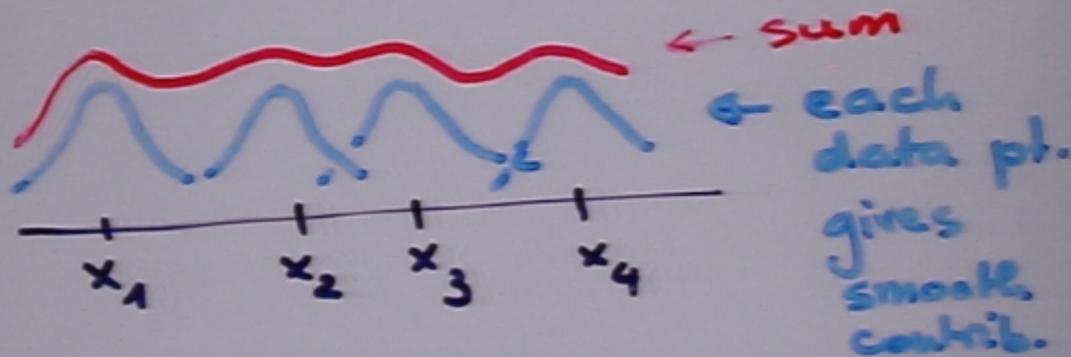


- adaptive grid

Frazer & Swinney
Darbellay & Vajda

still large systematic error
large statistical error
difficult in high dim's
... but fast!

- Kernel methods



very large systematic errors,
unless width of kernels
chosen optimally :

too wide \rightarrow poor resolution
too narrow \rightarrow entropies
underestimated

optimal choice ??

- estimators based on
nearest neighbour statistics

!!

too wide \rightarrow poor resolution
too narrow \rightarrow entropies
underestimated

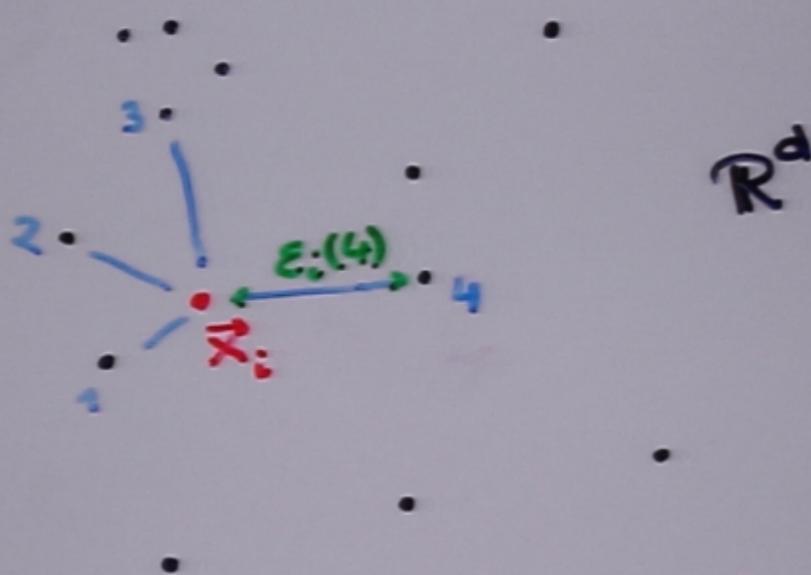
optimal choice ??

- estimators based on
nearest neighbour statistics

!!

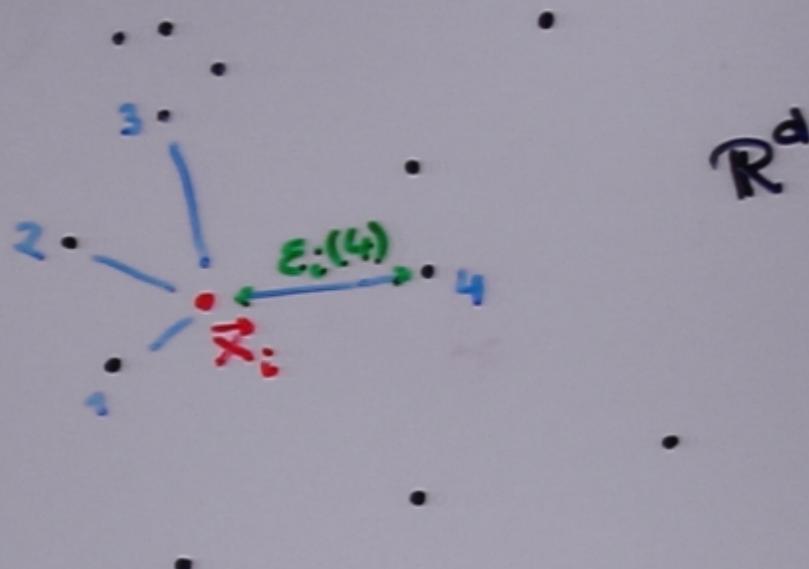
entropy:

$$H(X) = - \int dx \rho(x) \log \rho(x)$$



Kozachenko - Leonenko :

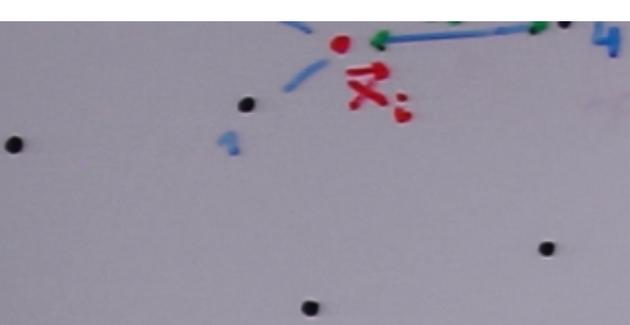
$$\hat{H}_k(x) = -\psi(k) + \psi(N) + \log c_d + \frac{d}{N} \sum \log \epsilon_i(k)$$



Kozachenko-Leonenko :

$$\hat{H}_k(x) = -\Psi(k) + \Psi(N) + \log c_d + \frac{d}{N} \sum_{|z_i| \leq k} \log \epsilon_i(k)$$

$$\Psi(x) = \frac{d}{dx} \ln \Gamma(x) \quad (\text{digamma-fct.})$$

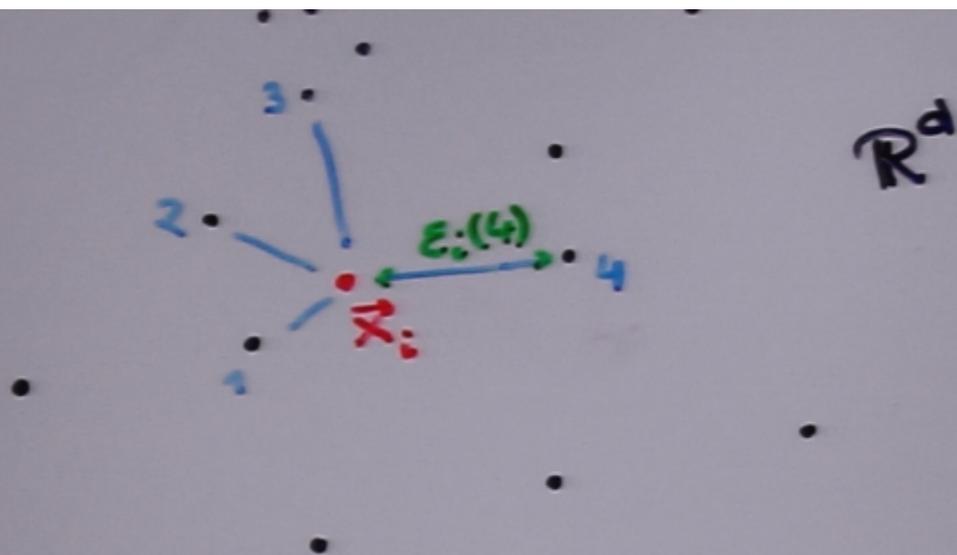


Kozachenko-Leonenko:

$$\hat{H}_k(x) = -\Psi(k) + \Psi(N) + \log c_d + \frac{d}{N} \sum_{i=1}^N \log \epsilon_i(k)$$

$$\Psi(x) = \frac{d}{dx} \ln \Gamma(x) \quad (\text{digamma-fct.})$$

c_d = volume of d -dim unit ball
 $\epsilon_i(k)$ = distance to k -th neighbor of x_i



Kozachenko-Leonenko:

$$\hat{H}_k(x) = -\Psi(k) + \Psi(N) + \log c_d + \frac{d}{N} \sum_{i=1}^N \log E_i(k)$$

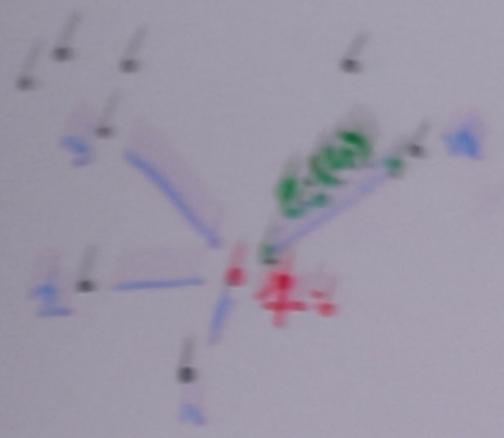
$$\Psi(x) = \frac{d}{dx} \ln \Gamma(x) \quad (\text{digamma-fct.})$$

c_d = volume of d -dim. unit ball
 $E_i(k)$ = distance to k -th neighbor of x_i

Handwritten notes:
 $H(x) = \dots$

$H(x) = \dots$

$P(x)$



Leonesko

$H(x) = \log \dots$

Estimation of $I(x, y)$:

- $\hat{I}_k(x, y) = \hat{H}_k(x) + \hat{H}_k(y) - \hat{H}_k(x, y)$

↑ bad !!

↑ $\epsilon(k)$ small

↑

↑ $\epsilon(k)$ large



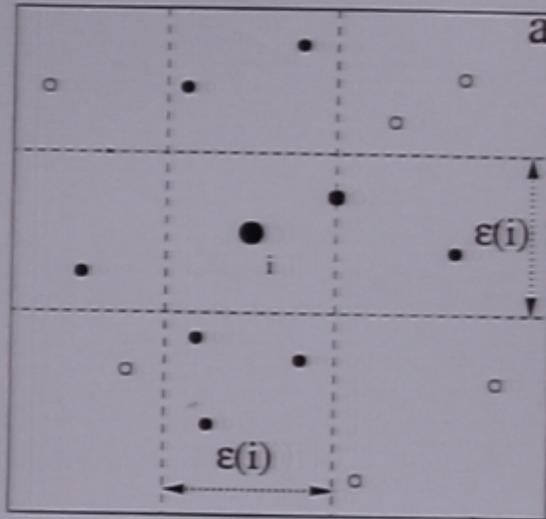
no cancelling
of errors

- better: use the same $\epsilon_i(k)$ for marginal & joint spaces !

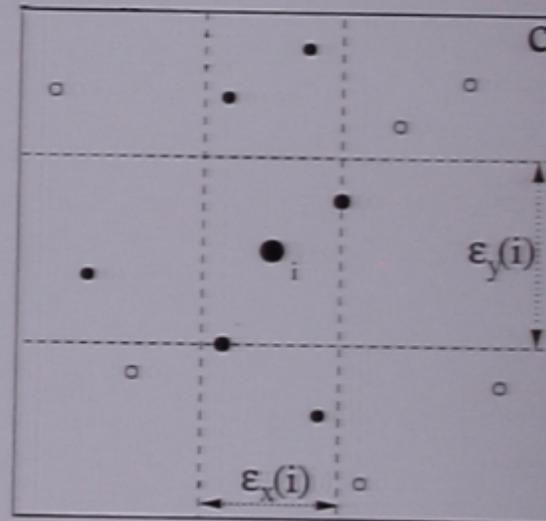
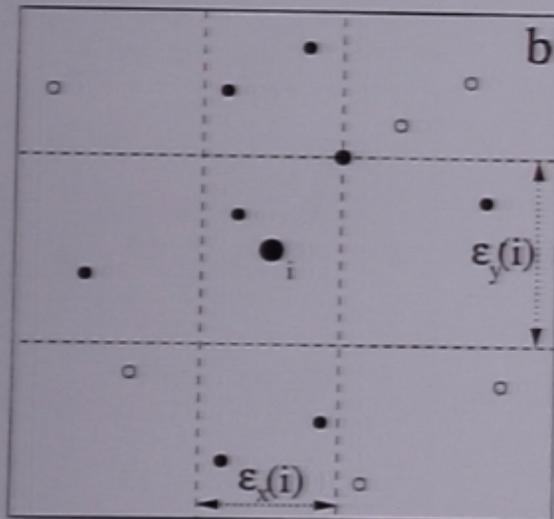
k : # neighbours in joint space

$n_x(i)$: - - - in X space

n_y - - - Y



hyper-cubic



rec

Handwritten mark, possibly initials or a signature, in blue ink.

• better: use the same $E_i(k)$
for marginal & joint
spaces !

k : # neighbours in joint space

$n_x(i)$: - - - in X space

n_y - - - Y

$$\Rightarrow \hat{I}(X, Y) = \psi(k) - \frac{1}{k} +$$
$$- \langle \psi(n_x) + \psi(n_y) \rangle$$
$$+ \psi(N)$$

no cancelling
of errors

- better: use the same $E_i(k)$
for marginal & joint
spaces !

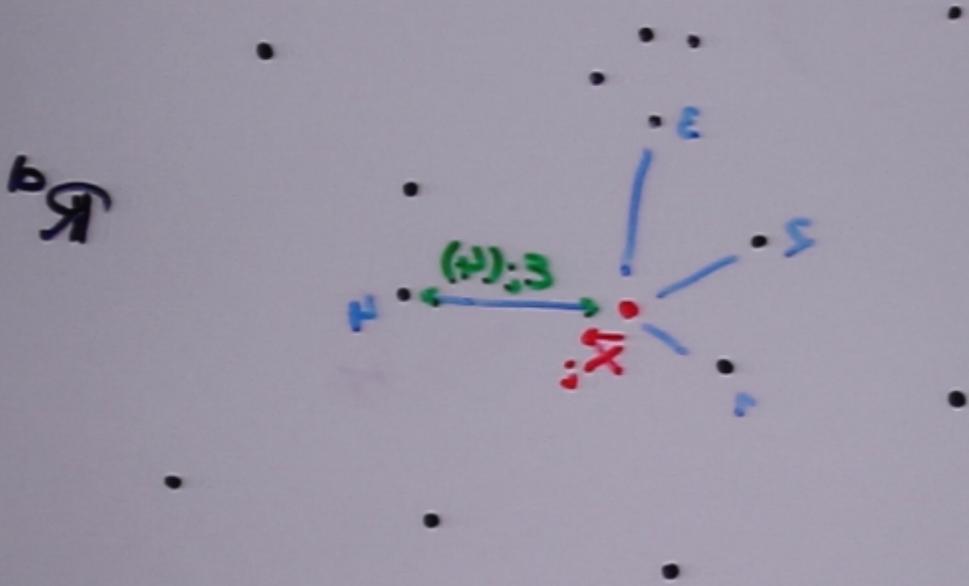
k : # neighbours in joint space

$n_x(i)$: - " - in X space

n_y - " - Y

$$\Rightarrow \hat{I}(x,y) = \psi(k) - \frac{1}{k} +$$
$$- \langle \psi(n_x) + \psi(n_y) \rangle$$
$$+ \psi(N)$$

$$\hat{H}(x) = \sum_{i=1}^n \psi_i(x) \hat{c}_i$$



Kozachenko-Leonenko:

$$\hat{H}_n(x) = \sum_{i=1}^n \psi_i(x) \hat{c}_i + \log c_n$$

no cancelling
of errors

- better: use the same $E_i(k)$
for marginal & joint
spaces !

k : # neighbours in joint space

$n_x(i)$: - " - in X space

n_y - " - Y

$$\Rightarrow \hat{I}(x,y) = \psi(k) - \frac{1}{k} + \\ - \langle \psi(n_x) + \psi(n_y) \rangle \\ + \psi(N)$$

no cancelling
of errors

- better: use the same $E_i(k)$
for marginal & joint
spaces !

k : # neighbours in joint space

$n_x(i)$: - - - in X space

n_y - - - Y

$$\Rightarrow \hat{I}(x,y) = \psi(k) - \frac{1}{k} +$$
$$- \langle \psi(n_x) + \psi(n_y) \rangle$$
$$+ \psi(N)$$

one uses (hyper-) cubes
(= balls in max norm)

other uses (hyper-) rectangles

other version slightly different

Both versions:

- small statistical errors
- fast, if efficient data structures used
- robust against outliers
- if degeneracies due to discretization
⇒ randomize by adding small noise

other version slightly different

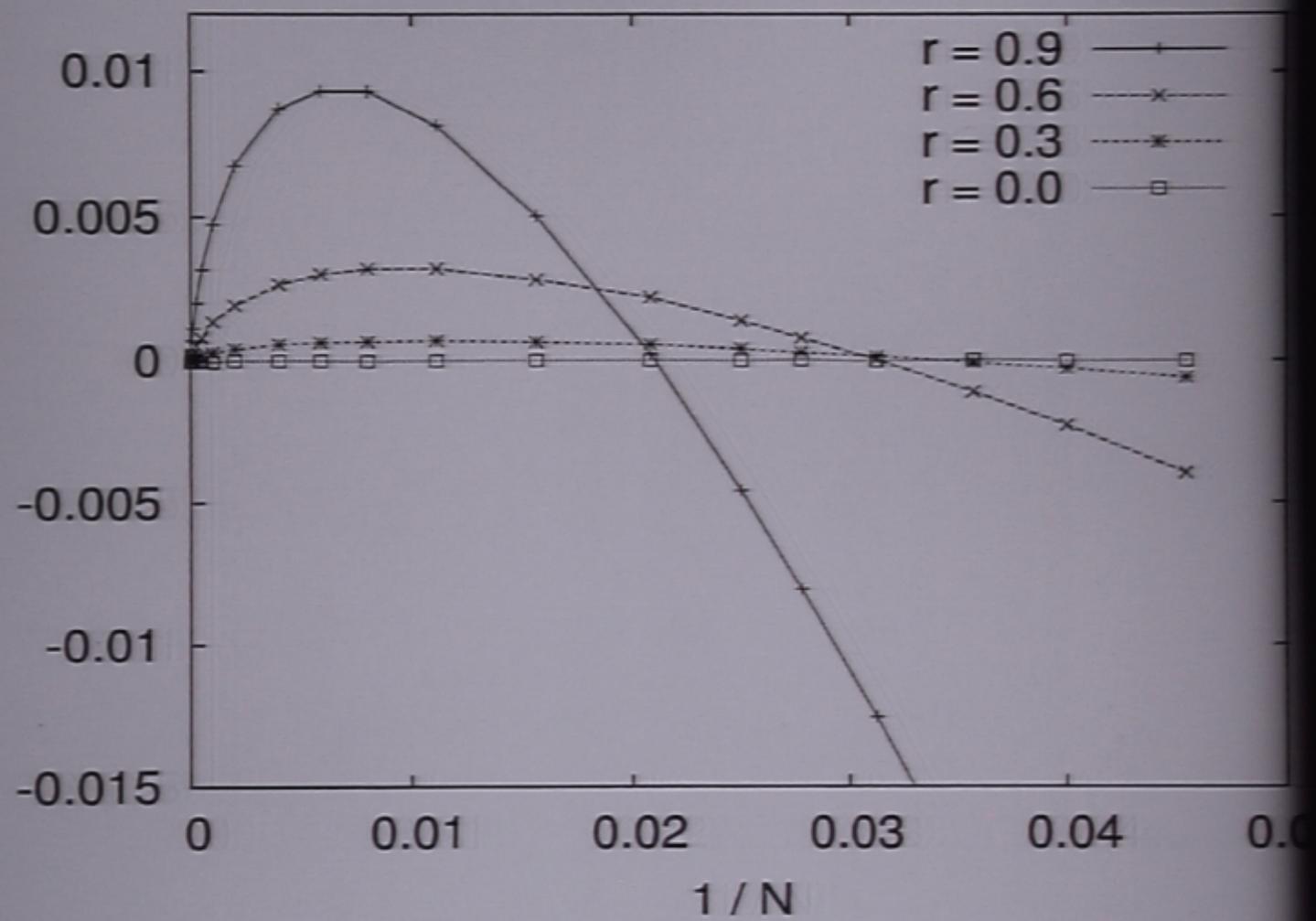
Both versions:

- small statistical errors
- fast, if efficient data structures used
- robust against outliers
- if degeneracies due to discretization
⇒ randomize by adding small noise
- small systematic errors
- **no** systematic errors, if X, Y

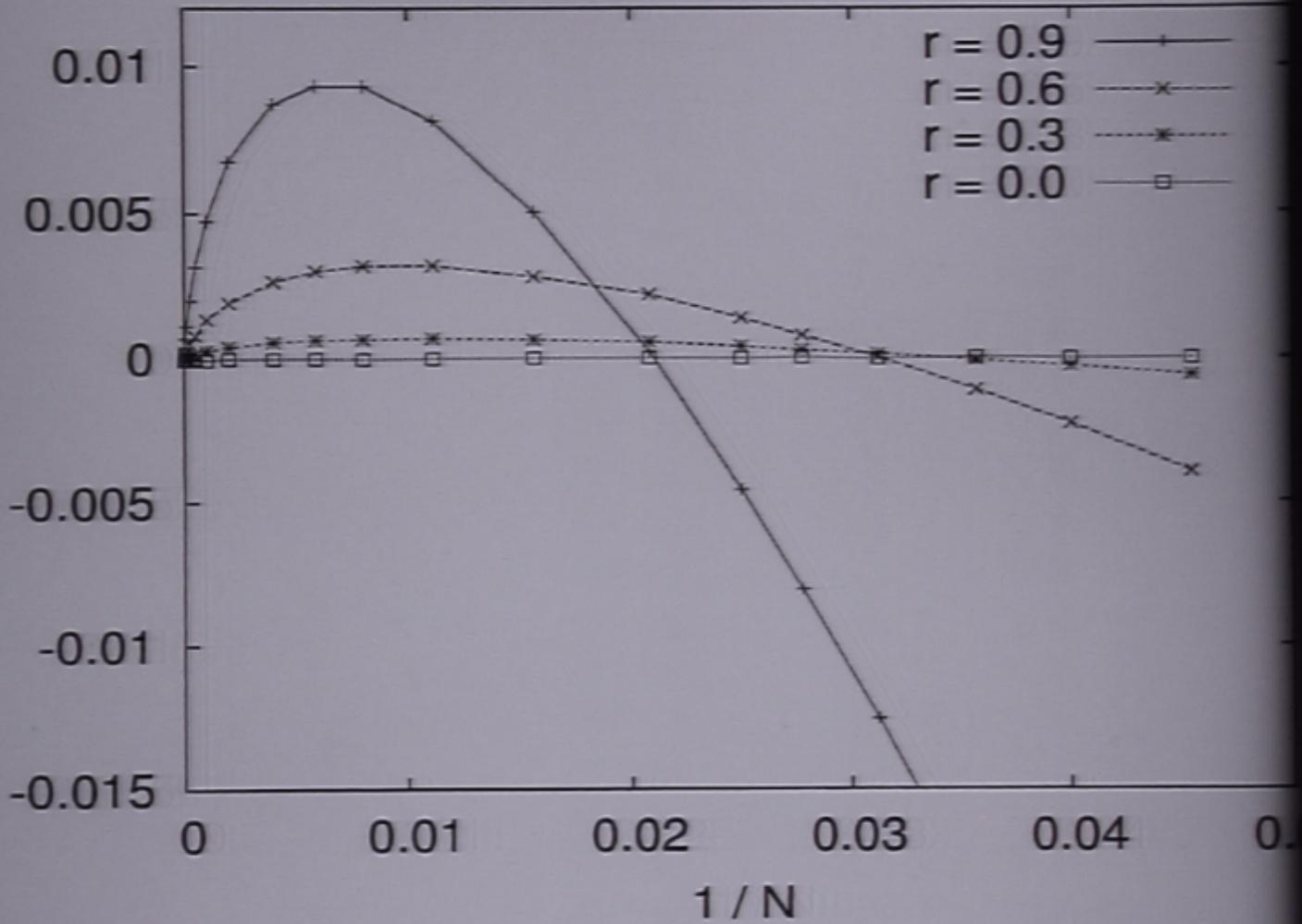
Both versions :

- small statistical errors
- fast, if efficient data structures used
- robust against outliers
- if degeneracies due to discretization
⇒ randomize by adding small noise
- small systematic errors
- no systematic errors, if X, Y independent (numerical)

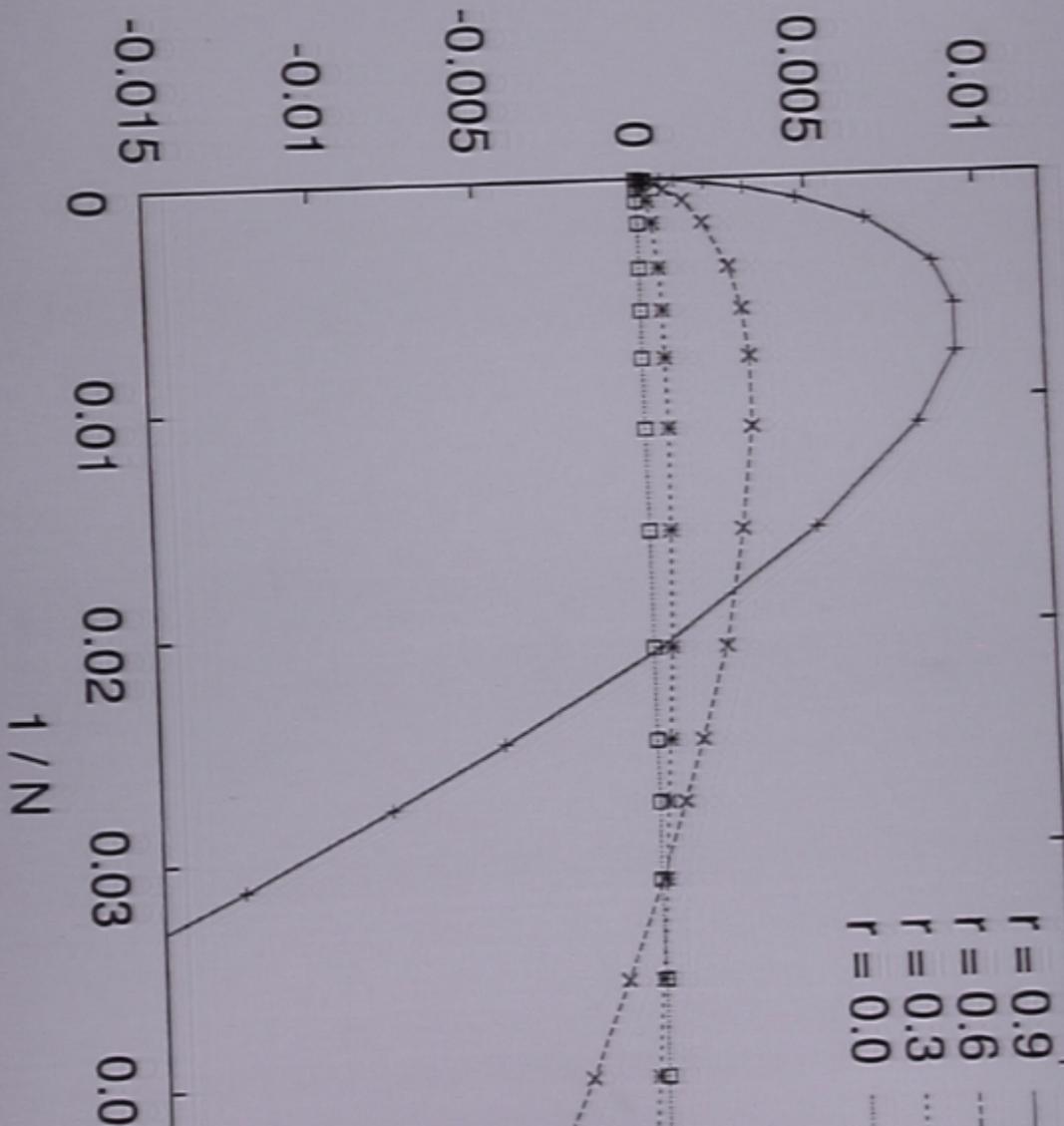
$$\langle I^{(2)}(X, Y) \rangle + \frac{1}{2} \log(1-r^2)$$



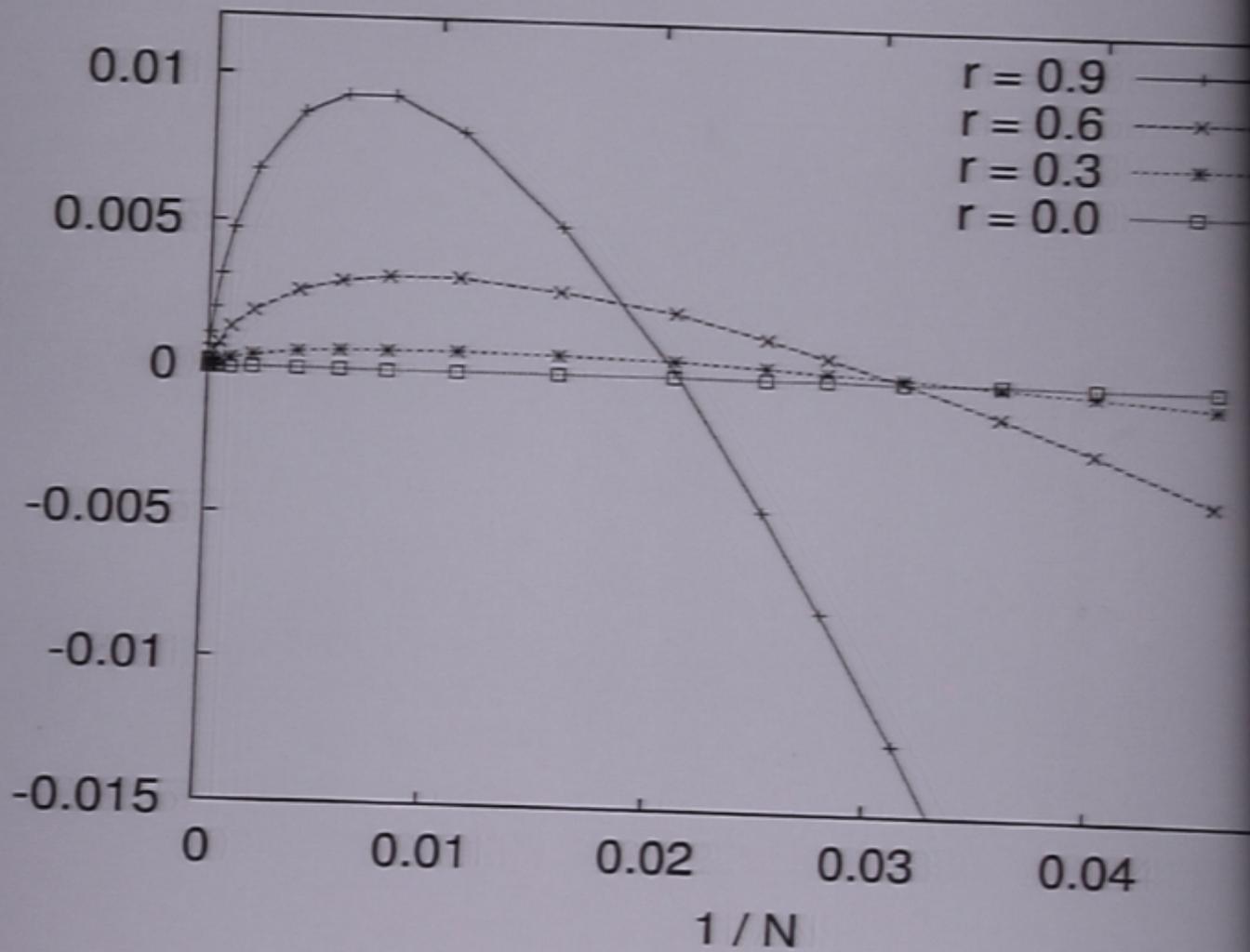
$$\langle I^{(2)}(X, Y) \rangle + \frac{1}{2} \log(1-r^2)$$



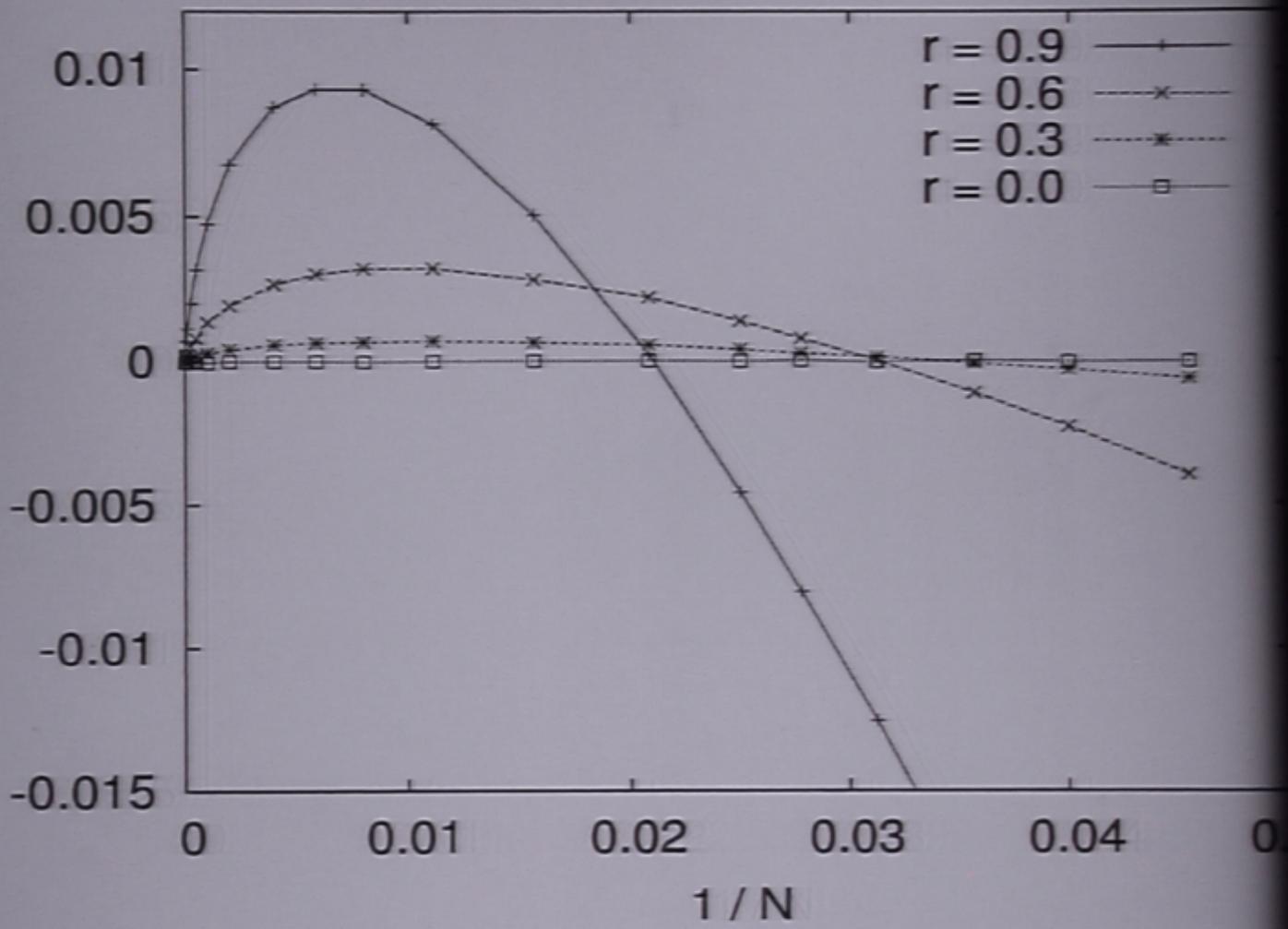
$$\langle l^{(2)}(X, Y) \rangle + 1/2 \log(1-r^2)$$

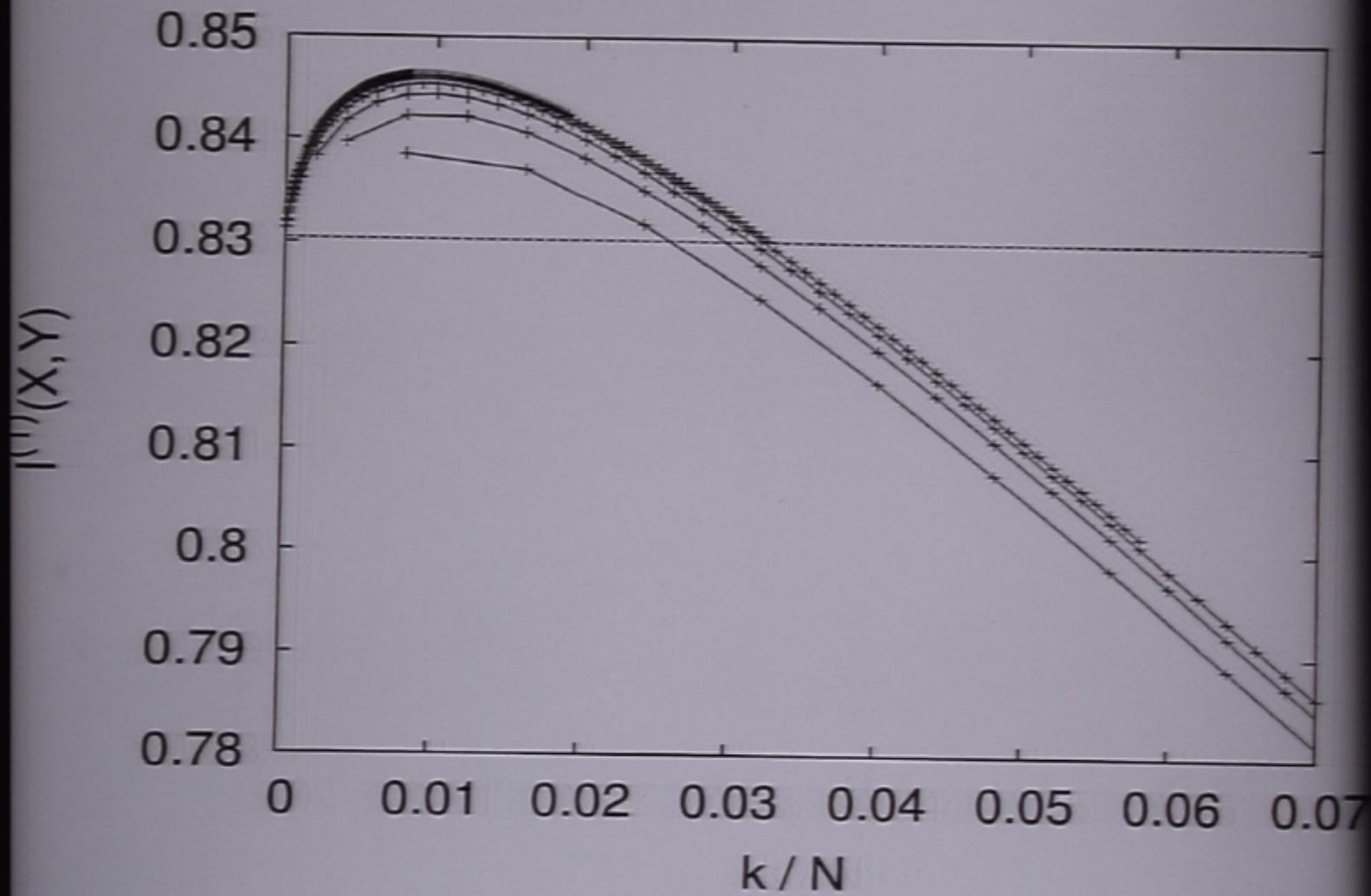


$$\langle I^{(2)}(X, Y) \rangle + \frac{1}{2} \log(1-r^2)$$



$$\langle I^{(2)}(X, Y) \rangle + \frac{1}{2} \log(1-r^2)$$

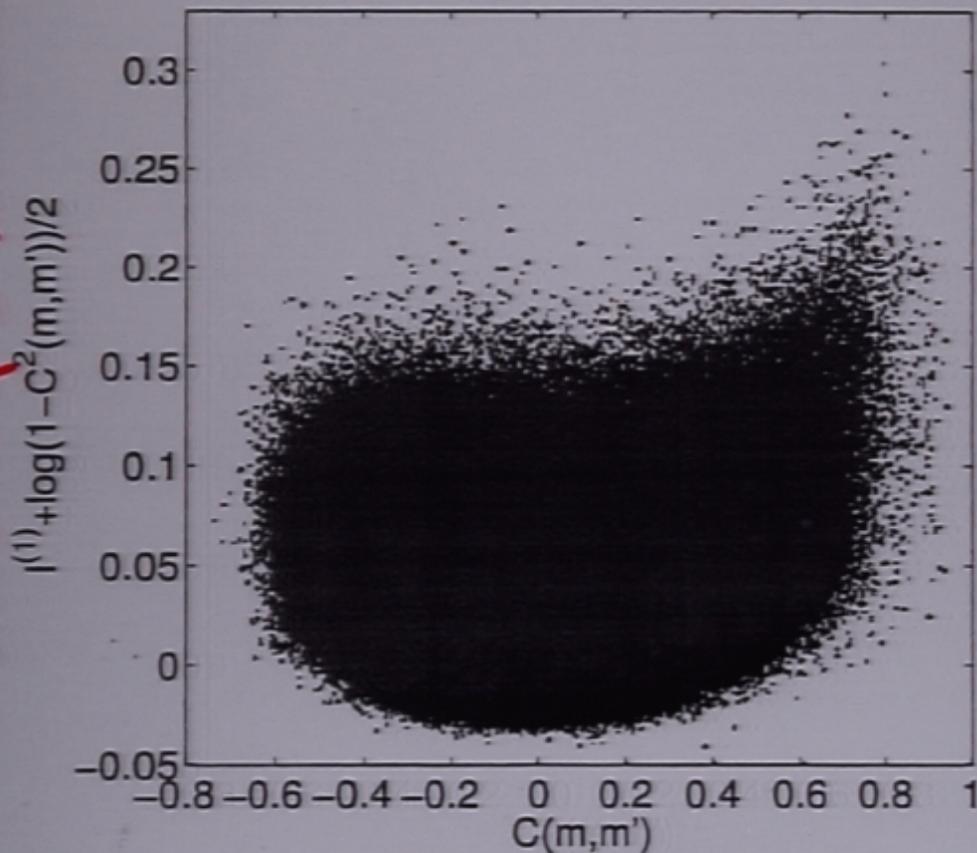




$m' = 1, \dots, 6000$ in yeast

dependencies measured via
expression ratios of 300 closely
related genomes (mostly a few mutations)

deviation of \bar{I} from lower bd.

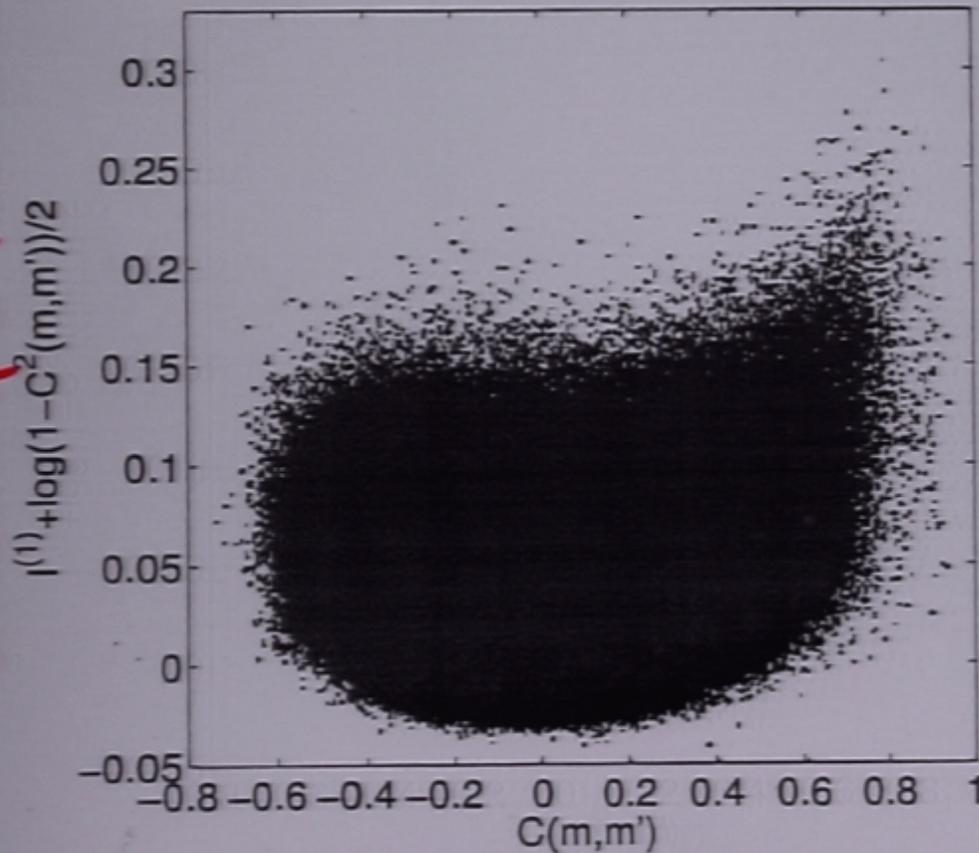


linear correl.

$m = 1, \dots, 6000$ } open reading frames
 $m' = 1, \dots, 6000$ } in yeast

dependencies measured via
expression ratios of 300 closely
related genomes (mostly a few mutations)

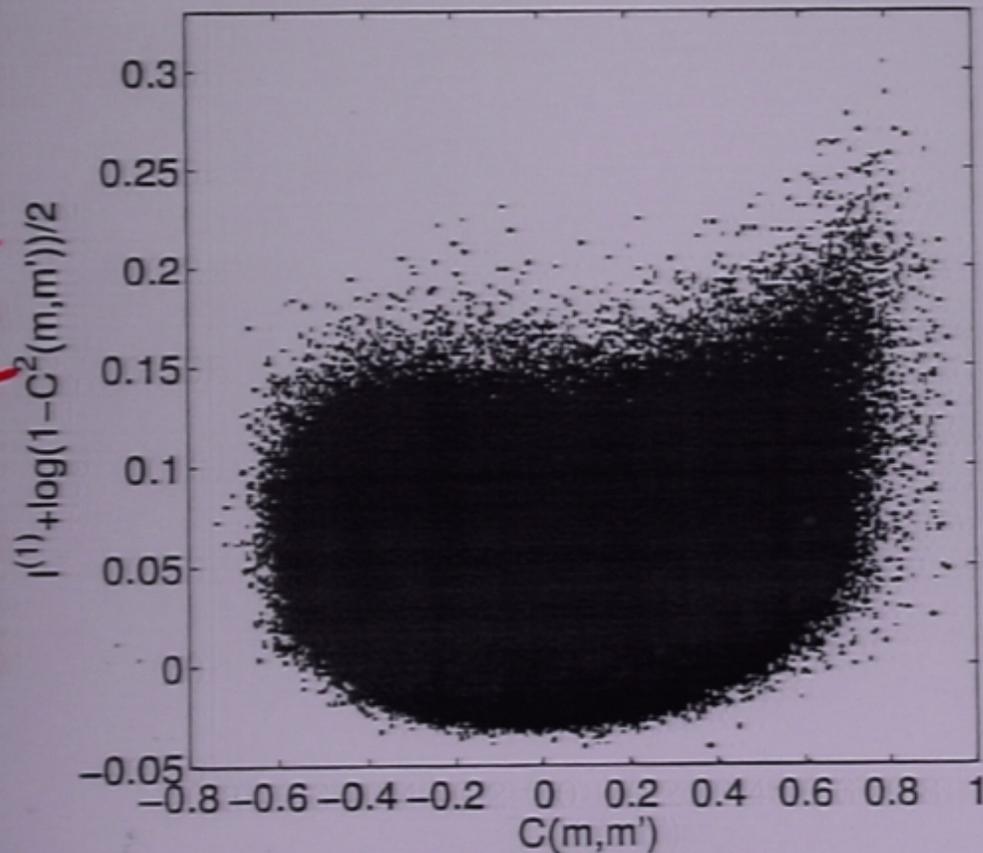
deviation of \hat{I} from lower bd.



$m = 1, \dots, 6000$ } open reading frames
 $m' = 1, \dots, 6000$ } in yeast

dependencies measured via
expression ratios of 300 closely
related genomes (mostly a few mutations)

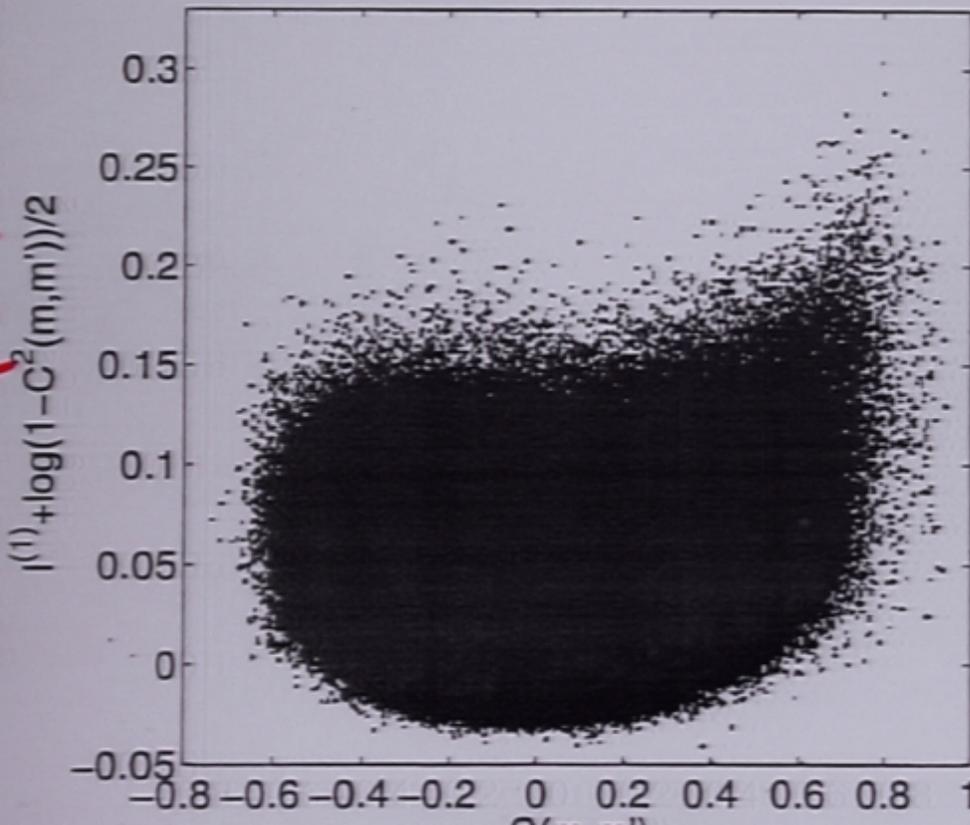
deviation of \hat{I} from lower bd.



$m = 1, \dots, 6000$ } open reading frames
 $m' = 1, \dots, 6000$ } in yeast

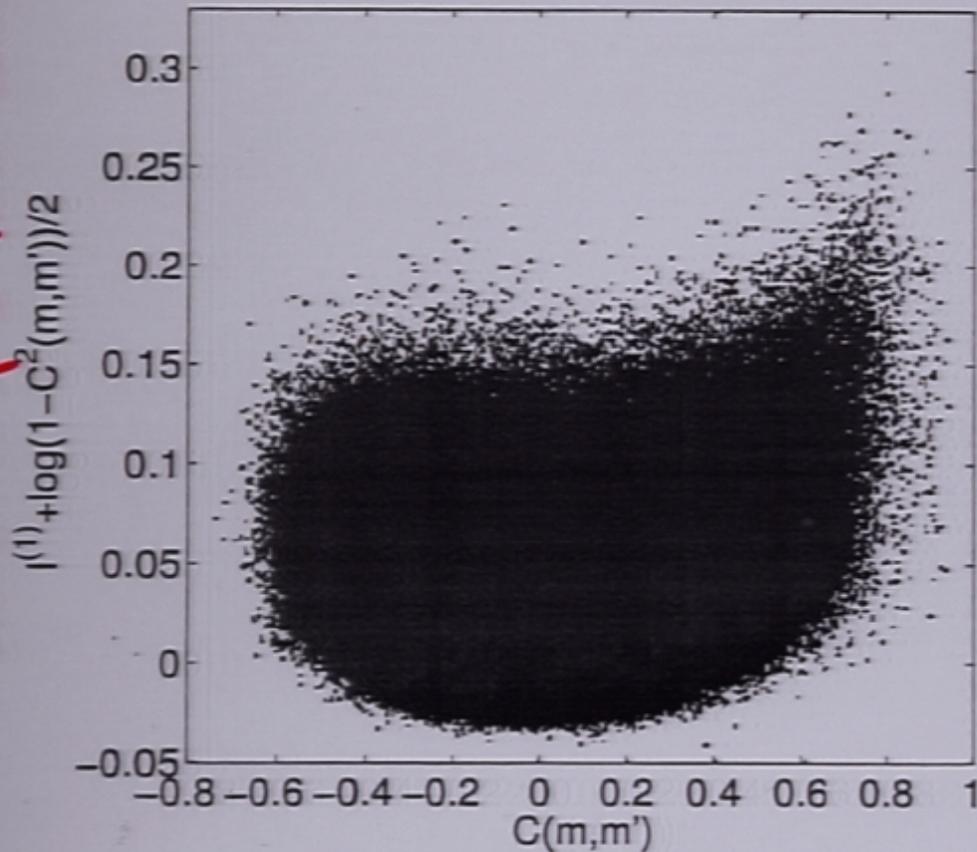
dependencies measured via
expression ratios of 300 closely
related genomes (mostly a few mutations)

deviation of \bar{I} from lower bd.



dependencies measured via
expression ratios of 300 closely
related genomes (mostly a few mutations)

deviation of \bar{I} from lower bd.



→
linear correl.
coeff.

Mutual Information based
Least dependent
Independent Component
Analysis

(MILCA)

Common problem in applied
Sciences:

measured data

= superposition of several
sources

Examples:

- signal + noise
- several speakers ("party problem")
- partial overlap between cellular phone freq. bands
- spectra emitted from chemical mixtures

Can sources be disentangled,
just from the observed data
alone ??

(no a priori knowledge of
how they were mixed,
no knowledge about number
of different sources, ...)

"blind source separation"

BSS

Only possible, if

Can sources be disentangled,
just from the observed data
alone ??

(no a priori knowledge of
how they were mixed,
no knowledge about number
of different sources, ...)

"blind source separation"

BSS

Only possible, if

no knowledge about number
of different sources, ...)

"blind source separation"

BSS

Only possible, if

- sources have different statistics
- or/and are independent

Simplest possibility:

- linear superposition
- statistically strictly independent sources
- # sources = # channels
- different "training tuples"
(i.e., signal vectors at diff. times)
are i.i.d

Simplest possibility:

- linear superposition
- statistically strictly independent sources
- # sources = # channels
- different "training tuples"
(i.e., signal vectors at diff. times)
are i.i.d

$$x_i(t) = \sum_{k=1}^n W_{ik} s_k(t) \quad i=1, \dots, n$$

↑
measure-
ments

↑
mixing
matrix
(instantan.,
stationary,
non-singular)

↑
sources
(i.i.d.,
independent)

$W = ?$
 $s = ?$

$$x_i(t) = \sum_{k=1}^n W_{ik} s_k(t) \quad i=1, \dots, n$$

↑
 measurements

mixing matrix
 (instantan.,
 stationary,
 non-singular)

↑
 sources
 (i.i.d.,
 independent)

$$W = ?$$

$$s = ?$$

obvious formal sol'n to simplest problem:

$$\vec{s}(t) = W^{-1} \vec{x}(t)$$

such that

$$s_i(t), s_k(t) \quad \text{independent} \\ \forall (i,k)$$

What does this mean?

E.g.

$$\langle s_i s_k \rangle - \langle s_i \rangle \langle s_k \rangle = 0$$

(uncorrelated :)

$$C_{ij}(x) \equiv \langle x_i x_j \rangle - \langle x_i \rangle \langle x_j \rangle$$

$$\vec{s}(t) = W^{-1} \vec{x}(t)$$

Such that

$$s_i(t), s_k(t) \quad \text{independent} \\ \neq (i, k)$$

What does this mean?

E.g. $\langle s_i s_k \rangle - \langle s_i \rangle \langle s_k \rangle = 0$
(uncorrelated)

$$C_{ij}(x) \equiv \langle x_i x_j \rangle - \langle x_i \rangle \langle x_j \rangle$$

$$W^{-1} C(x) W \stackrel{!}{=} \text{diagonal}$$

\Rightarrow PCA

If only covariance matrix is used as dependency measure \Rightarrow PCA

BUT: uncorrelated \neq independent!

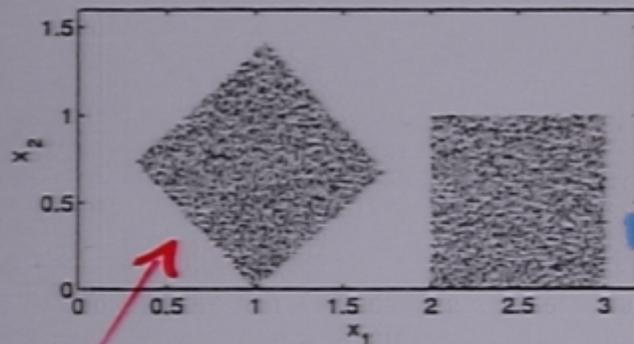


Abbildung 3.1: Zwei gleichförmig verteilte Zufallsvariable, wobei sich die linke von der rechten Abbildung nur durch eine Rotation unterscheidet. Beide Bilder stellen unkorrelierte Variable dar, aber nur rechts sind sie auch unabhängig.

x, y are uncorrelated

x, y are statistically independent

ICA:

assumes same linear mixing ansatz as PCA, but

goes beyond 2nd order statistics
for finding W

PCA does not give unique $A = W^{-1}$:

- permutation of indices
- re-scaling of axes

Use A such that $C(s) = \text{const} \cdot \mathbf{1}$

\Rightarrow "pre-whitening"

assumes same linear mixing
ansatz as PCA, but

goes beyond 2nd order statistics

for finding W

PCA does not give unique $A = W^{-1}$:

- permutation of indices
- re-scaling of axes

Use A such that $C(s) = \text{const} \times \mathbb{1}$

\Rightarrow "pre-whitening"

ansatz as PCA, but

goes beyond 2nd order statistics

for finding W

PCA does not give unique $A = W^{-1}$:

- permutation of indices
- re-scaling of axes

Use A such that $C(s) = \text{const} = 1$

⇒ "pre-whitening"

- ⑥ after pre-whitening, there is still an unspecified

goes beyond 2nd order statistics
for finding W

PCA does not give unique $A \neq W^{-1}$:

- permutation of indices
- re-scaling of axes

Use A such that $C(s) = \text{const} \times \mathbf{1}$

\Rightarrow "pre-whitening"

- after pre-whitening, there is still an unspecified rotation $R \in O(n)$ free

ICA proper specifies R

up to

- permutation of sources
- common re-scaling of all sources

Jutten & Herault ~ 1980

Bell & Sejnowski ~ 1990

Cardoso ≥ 1990

⋮

A. Hyvärinen, J. Karhunen, E. Oja
"Independent Component Anal."
Wiley

- permutation of sources
- common re-scaling of all sources

Jutten & Herault ~ 1980
Bell & Sejnowski ~ 1990
Cardoso \geq 1990
⋮

A. Hyvärinen, J. Karhunen, E. Oja
"Independent Component Anal."
Wiley 2001

many public domain & proprietary
classics.

Basic Milca:

minimize $\hat{I}(s_1, \dots, s_n)$ by rotations
in subspaces (after whitening)

choose in each subspace optimal
angle by fitting \hat{I} with short
Fourier series

• Significance test:

check variation of $\hat{I}(s_1, \dots, s_n)$
with angles,

to find angles which cannot be
given unambiguously
(no unique BSS possible)

- Basic Milca:

minimize $\hat{I}(s_1 \dots s_n)$ by rotations
in subspaces (after whitening)

choose in each subspace optimal
angle by fitting \hat{I} with short
Fourier series

- Significance test:

check variation of $\hat{I}(s_1 \dots s_n)$
with angles,

to find angles which cannot be
given unambiguously
(no unique BSS possible)

- Independence test :

check whether final
 $\hat{I}(s_i; s_j) \cong 0 \quad \forall i, j$

- Multi-variate MILCA

if \exists cluster $\{s_1 \dots s_n\}$

of dependent sources,

then consider this as one

k-variate source

Use grouping property for a
complete clustering of all sources
("mutual information based
clustering", MIC :

$\exists \rightarrow$ cluster $\{S_1 \dots S_n\}$
of dependent sources,
then consider this as one
k-variate source

Use grouping property for a
complete clustering of all sources
("mutual information based
clustering", MIC :
A. Kraskov et al)

reconstruct contributions of
each relevant cluster to every
channel

- MILCA with instantaneous mixing
time correlated signals

$x_i(t), x_i(t')$ not indep.

$$\vec{S}(t) = A \vec{x}(t) \quad \text{as before} \\ \text{(no delay)}$$

Def:

$$\underline{s}_i(t) = (s_i(t-d\tau) \dots s_i(t))$$

$d+1 = d_{i, \dots}$
delay vectors

$$\left| \hat{I}(\underline{s}_1(t) \dots \underline{s}_n(t)) \stackrel{!}{=} \min \right|$$

minimize also dependence
between time-shifted
sources!

- MILCA with delayed mixing
& delayed \hat{I}

our best algorithm!

$$\hat{s}_i(t) = \sum_{j=1}^n \sum_{k=0}^d a_{ij}^k x_j(t - k\tau)$$

$i=1 \dots nd$

$$\hat{I}(\hat{s}_1 \dots \hat{s}_{nd}) \stackrel{!}{=} \min$$

ECG of pregnant woman

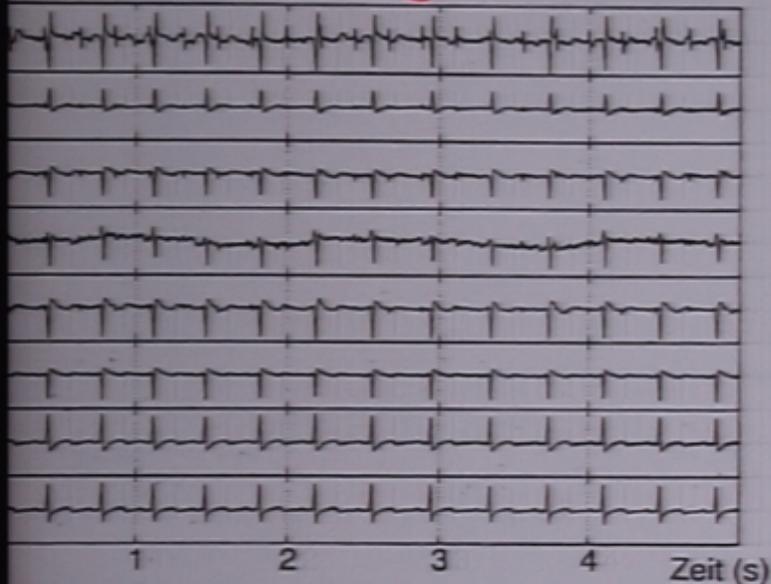


Abbildung 4.1: 8-Kanal EKG einer schwangeren Frau.

basic MILCA

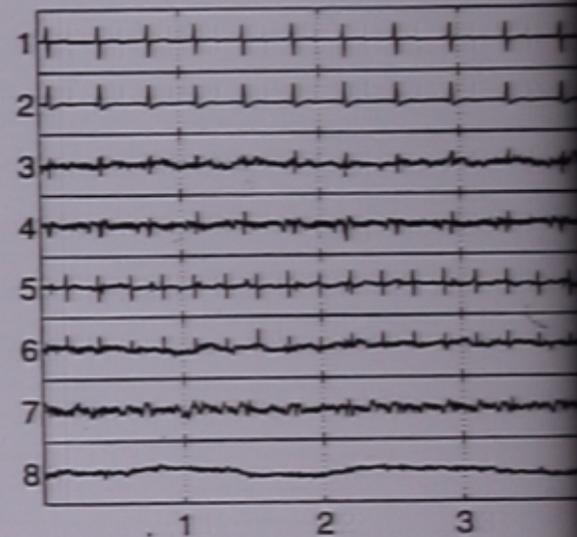
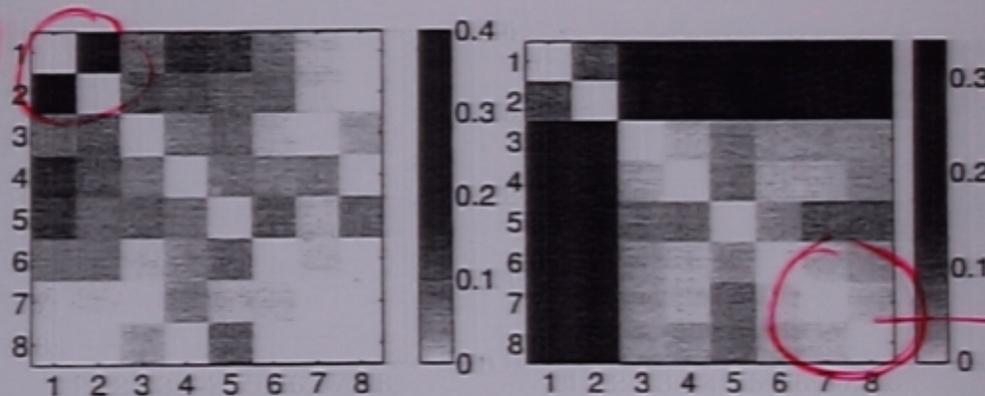


Abbildung 4.2: MILCA-Ausgabe: Komponenten nach der MILCA für die EKG-Signale aus Abb.(4.1).

strongly
correl.
cluster



no u
dece
of u

Abbildung 4.3: Linkes Bild: \hat{r}_{ij} zwischen allen paarweisen Kombinationen des Signals aus Abb.(4.2). Rechtes Bild: Quadratwurzel der Variabilität σ_{ij} von $\hat{r}_{ij}(\phi)$. In beiden Bildern ist die Diagonale auf Null gesetzt.

ECG of pregnant woman

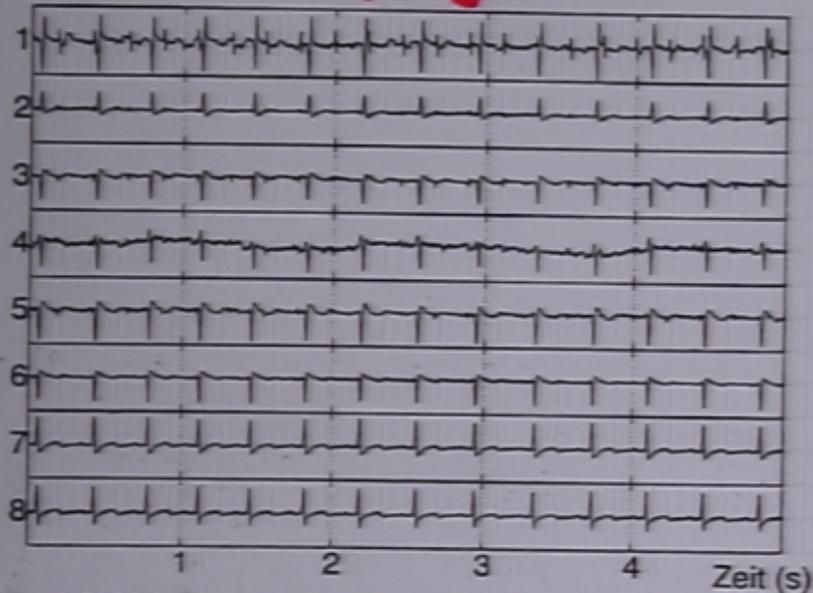


Abbildung 4.1: 8-Kanal EKG einer schwangeren Frau.

basic

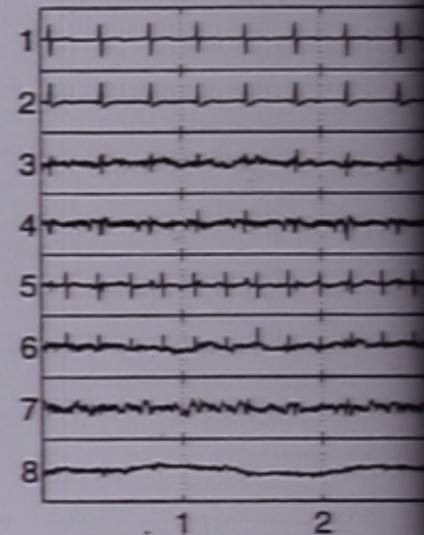


Abbildung 4.2: MILCA-Ausgabe: Komponenten für die EKG-Signale aus Abb.(4.1).

strongly
correl.
clusters

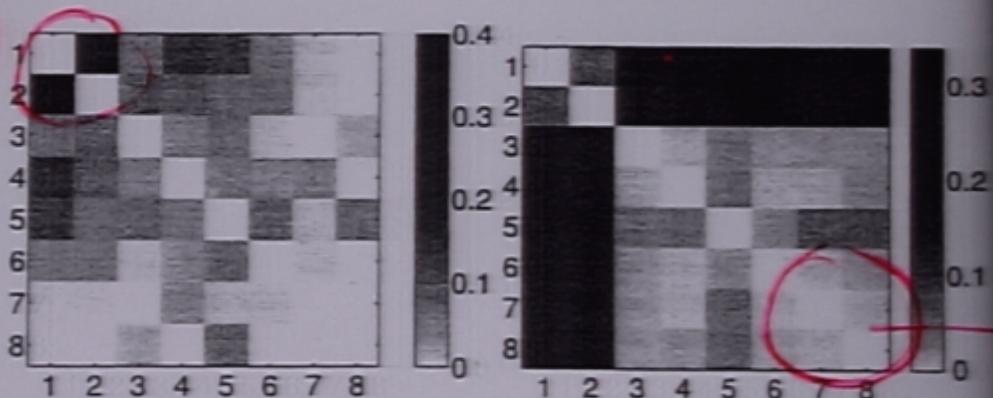
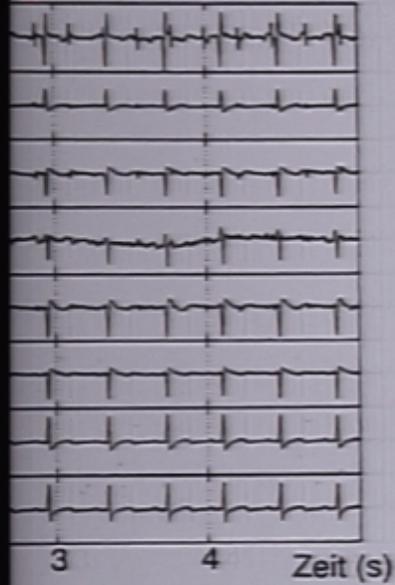


Abbildung 4.3: Linkes Bild: \bar{I} zwischen allen paarweisen Kombinationen des Signals aus Abb.(4.2). Rechtes Bild: Quadratwurzel der Variabilität σ_{ij} von $\bar{I}_{ij}(\varphi)$. In beiden Bildern ist die Diagonale auf Null gesetzt.

gnant woman



einer schwangeren Frau.

basic MILCA

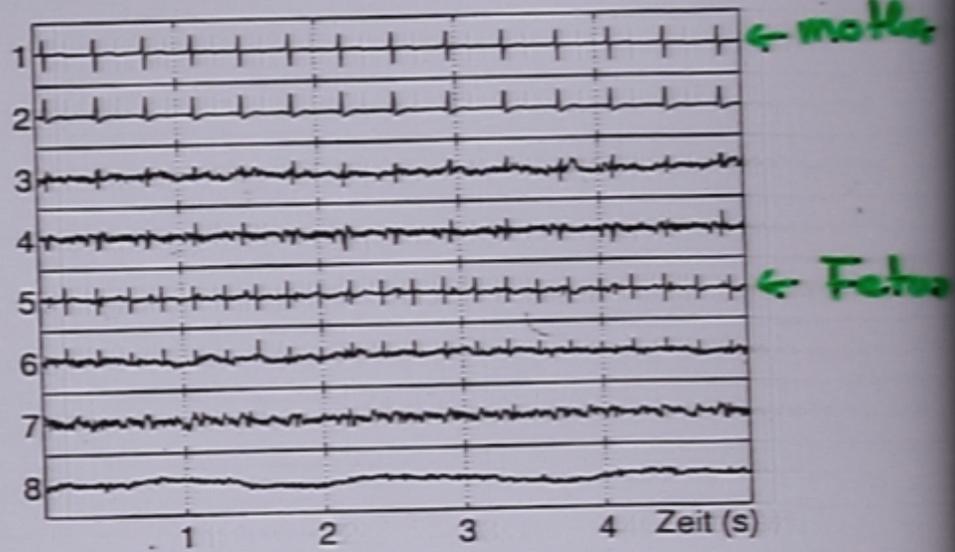
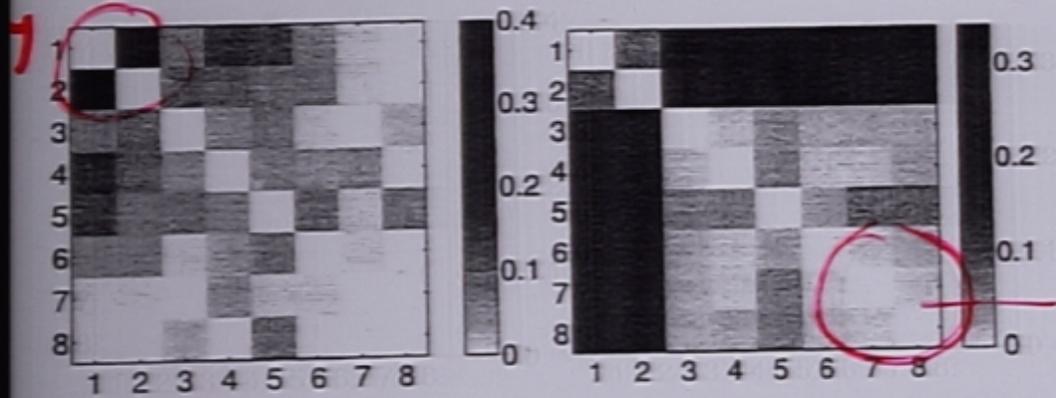


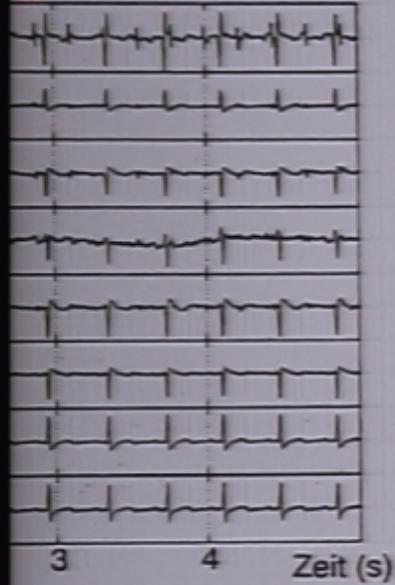
Abbildung 4.2: MILCA-Ausgabe: Komponenten nach der Minimierung von $I(X_1 \dots X_8)$ für die EKG-Signale aus Abb.(4.1).



no unique decomposition of noise sources

Abbildung 4.3: Linkes Bild: \hat{I} zwischen allen paarweisen Kombinationen des Signals aus Abb.(4.2). Rechtes Bild: Quadratwurzel der Variabilität σ_{ij} von $\hat{I}_{ij}(\phi)$. In beiden Bildern ist die Diagonale auf Null gesetzt.

gnant woman



einer schwangeren Frau.

basic MILCA

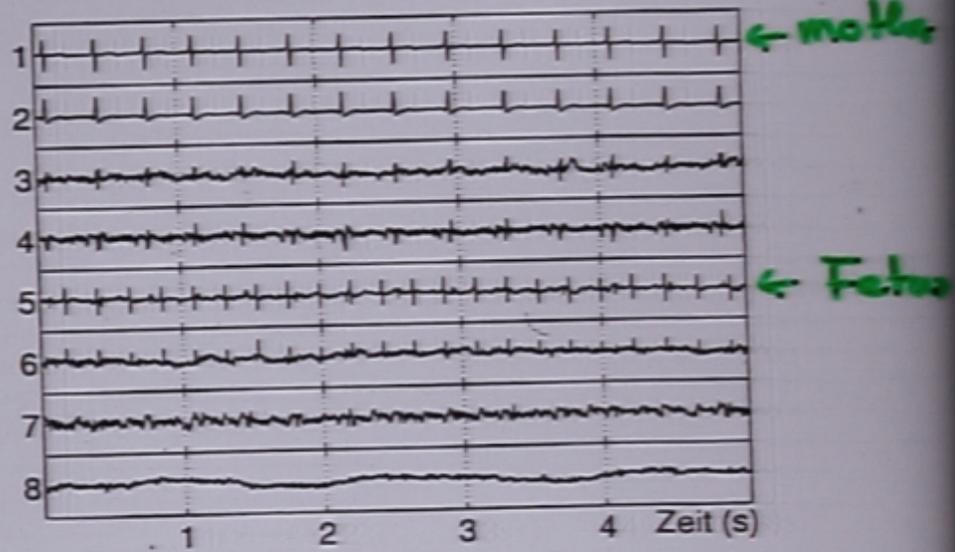
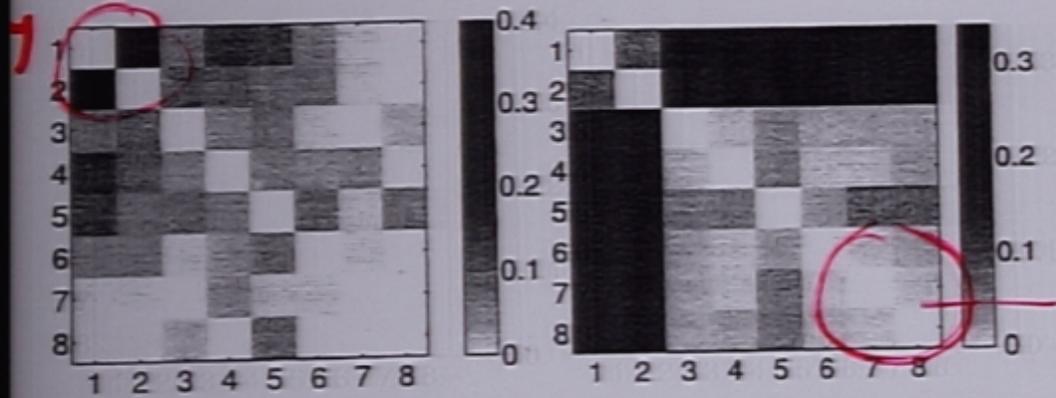


Abbildung 4.2: MILCA-Ausgabe: Komponenten nach der Minimierung von $I(X_1 \dots X_8)$ für die EKG-Signale aus Abb.(4.1).



no unique decomposition of noise sources

Abbildung 4.3: Linkes Bild: \hat{I} zwischen allen paarweisen Kombinationen des Signals aus Abb.(4.2). Rechtes Bild: Quadratwurzel der Variabilität σ_{ij} von $\hat{I}_{ij}(\phi)$. In beiden Bildern ist die Diagonale auf Null gesetzt.

After delay embedding (3 times) :
 8 channels \rightarrow 24 channels

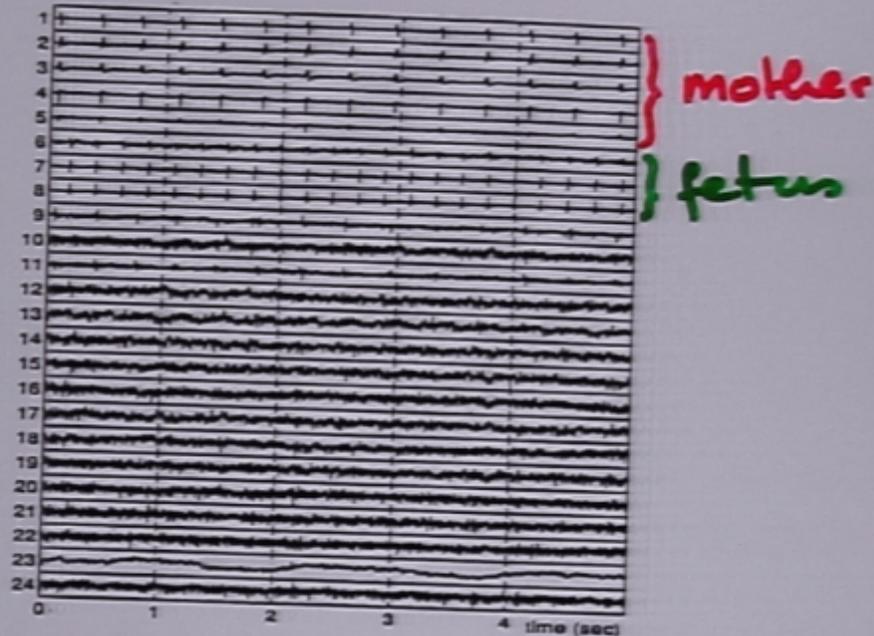


FIG. 17: MILCA output from the embedded eight channel ECG ($k=100, m=3$)

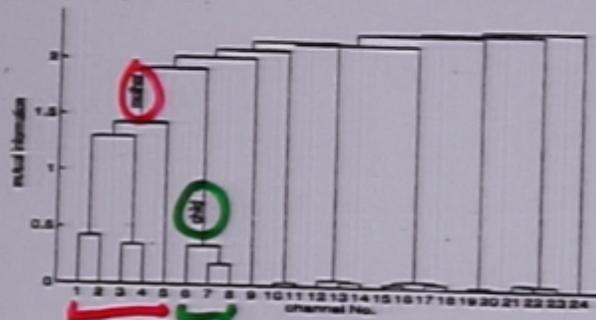


FIG. 18: Dendrogram for Fig. 17. Heights of each cluster correspond to $I(X_i, X_j)$ of the cluster ij ($k=3$).

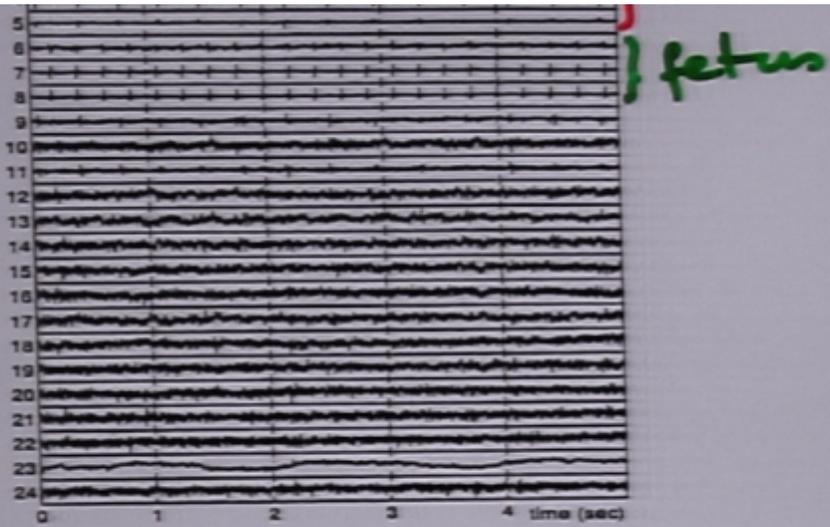


FIG. 17: MILCA output from the embedded eight channel ECG ($k=100, m=3$)

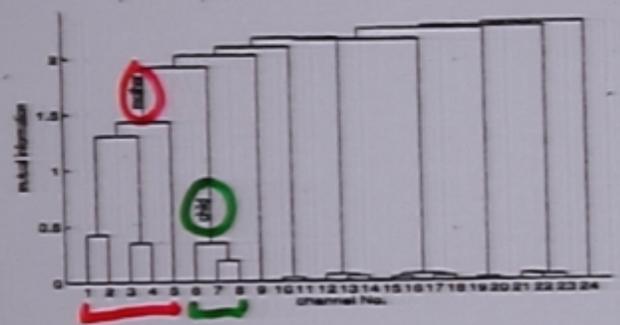


FIG. 18: Dendrogram for Fig. 17. Heights of each cluster correspond to $I(X_i, X_j)$ of the cluster ij ($k=3$).

8 channels → 24

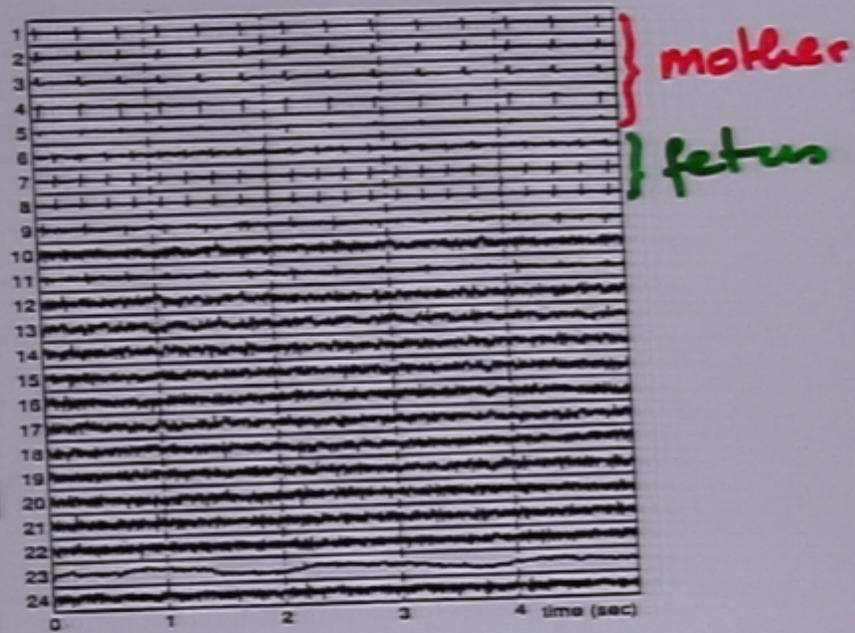


FIG. 17: MILCA output from the embedded eight channel ECG ($k=100, m=3$)

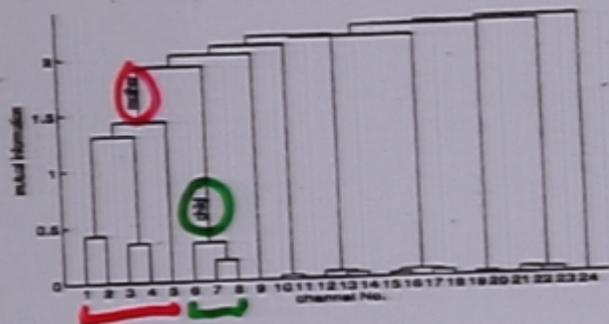


FIG. 18: Dendrogram for Fig. 17. Heights of each cluster correspond to $I(X_i, X_j)$ of the cluster ij ($k=3$).

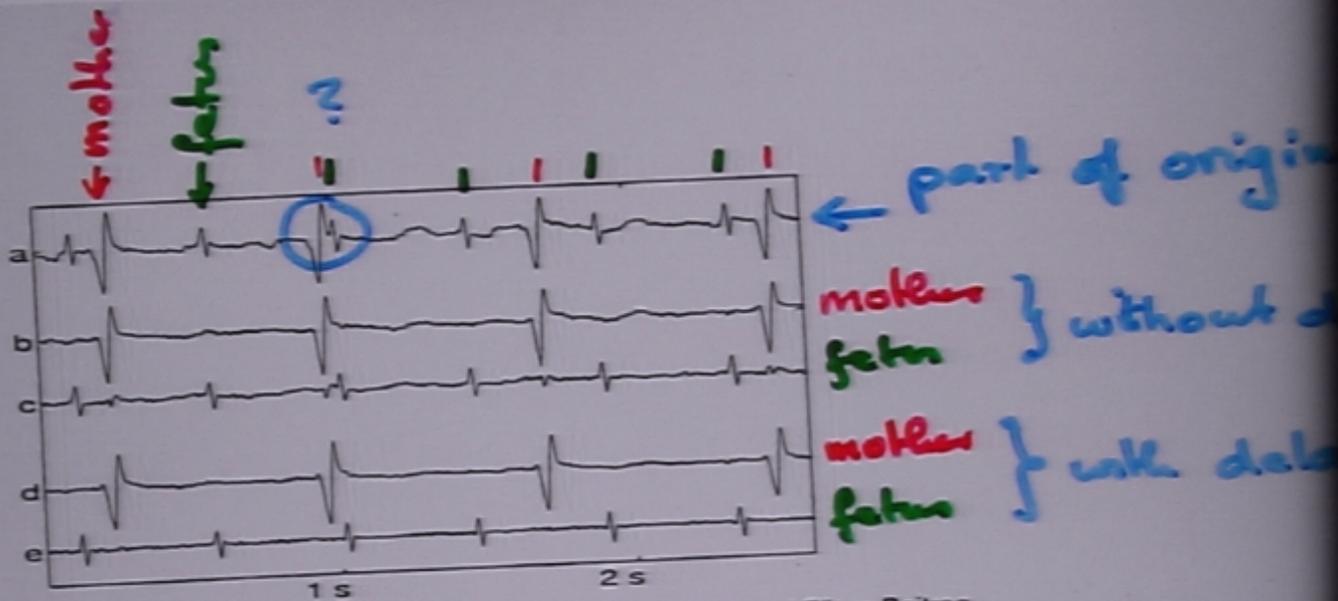


Abbildung 4.10: (a) Ausschnitt des originalen EKGs; (b,c) der Mutter- und Fötus-Beitrag, erhalten ohne Verzögerungseinbettung; (d,e) beide Beiträge bei Verwendung der Verzögerungseinbettung.

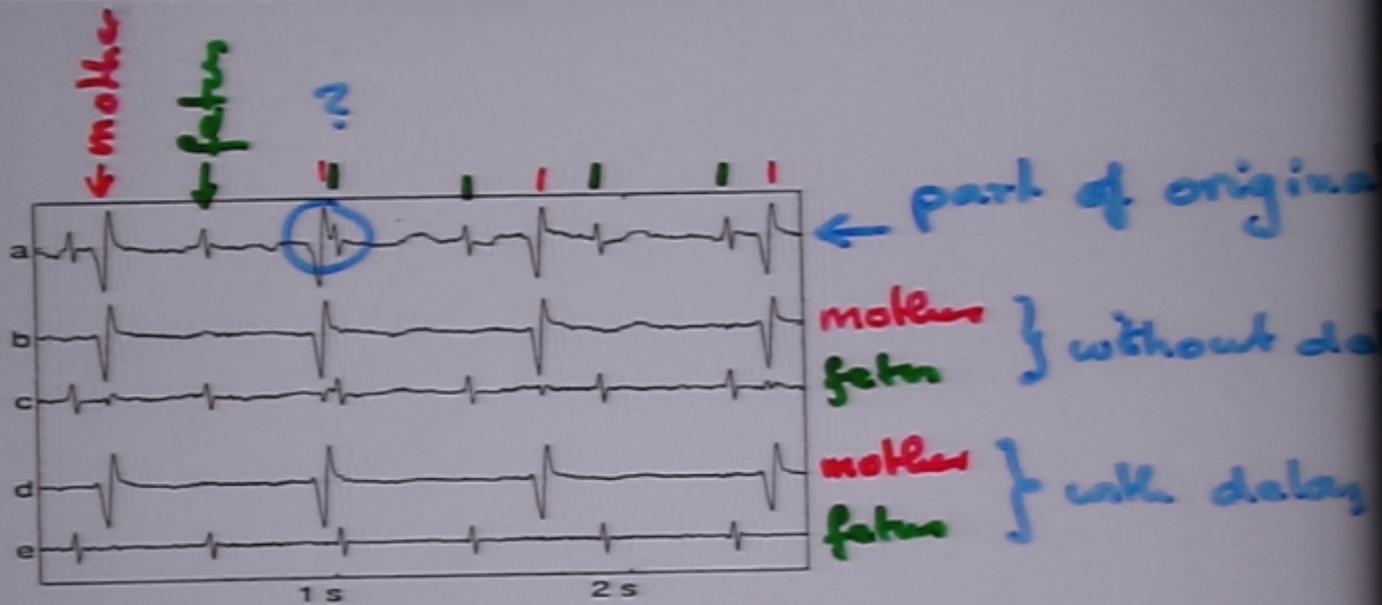


Abbildung 4.10: (a) Ausschnitt des originalen EKGs; (b,c) der Mutter- und Fötus-Beitrag, erhalten ohne Verzögerungseinbettung; (d,e) beide Beiträge bei Verwendung der Verzögerungseinbettung.