

Blind Source Separation:
PCA & ICA

HST.582J/6.555J/16.456J
Gari D. Clifford



gari [at] mit . edu
<http://www.mit.edu/~gari>





**Harvard-MIT
Health Sciences & Technology**

© G. D. Clifford 2005-2009

What is BSS?

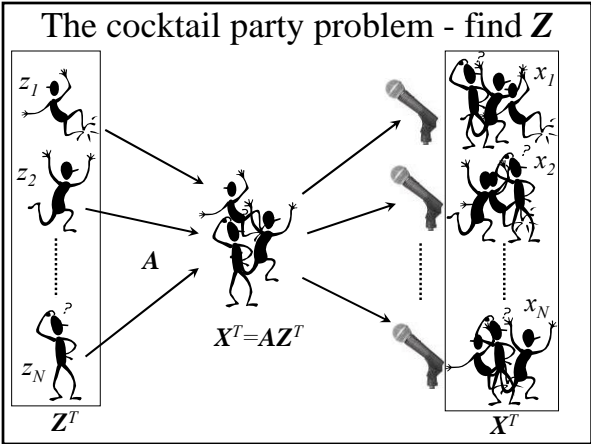
Assume an observation (signal) is a **linear** mix of >1 unknown **independent source** signals

The *mixing* (not the signals) is **stationary**

We have as many observations as unknown **sources**

To find sources in observations
- need to define a suitable **measure of independence**

... For example - the *cocktail party problem* (sources are speakers and background noise):

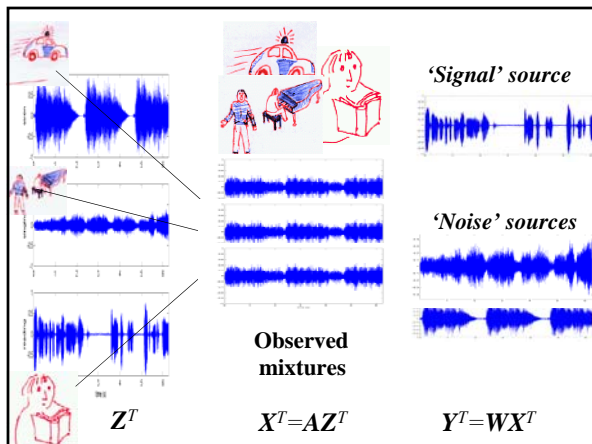


Formal statement of problem

- N independent sources ... \mathbf{Z}_{mn} ($M \times N$)
- linear square mixing ... \mathbf{A}_{mn} ($N \times N$)
(#sources=#sensors)
- produces a set of observations ... \mathbf{X}_{mn} ($M \times N$)
..... $\mathbf{X}^T = \mathbf{A}\mathbf{Z}^T$

Formal statement of solution

- 'demix' observations ... \mathbf{X}^T ($N \times M$)
into $\mathbf{Y}^T = \mathbf{W}\mathbf{X}^T$
 \mathbf{Y}^T ($N \times M$) $\approx \mathbf{Z}^T$ \mathbf{W} ($N \times N$) $\approx \mathbf{A}^{-1}$
- How do we recover the independent sources?
(We are trying to estimate $\mathbf{W} \approx \mathbf{A}^{-1}$)
- We require a **measure of independence!**



$$\begin{bmatrix} x_{11} & x_{12} & \dots & x_{1N} \\ x_{21} & x_{22} & \dots & x_{2N} \\ \vdots & \vdots & \ddots & \vdots \\ \vdots & \vdots & \ddots & \vdots \\ x_{M1} & x_{M2} & \dots & x_{MN} \end{bmatrix}^T = \begin{bmatrix} a_{11} & \dots & a_{N1} \\ \vdots & \ddots & \vdots \\ a_{1N} & \dots & a_{NN} \end{bmatrix} \begin{bmatrix} z_{11} & z_{12} & \dots & z_{1N} \\ z_{21} & z_{22} & \dots & z_{2N} \\ \vdots & \vdots & \ddots & \vdots \\ z_{M1} & z_{M2} & \dots & z_{MN} \end{bmatrix}^T$$

$$\mathbf{X}^T = \mathbf{A} \mathbf{Z}^T \quad \mathbf{Y} = \hat{\mathbf{Z}} \\ \mathbf{W} = \hat{\mathbf{A}}^{-1}$$

$$\begin{bmatrix} y_{11} & y_{12} & \dots & y_{1N} \\ y_{21} & y_{22} & \dots & y_{2N} \\ \vdots & \vdots & \ddots & \vdots \\ \vdots & \vdots & \ddots & \vdots \\ y_{M1} & y_{M2} & \dots & y_{MN} \end{bmatrix}^T = \begin{bmatrix} w_{11} & \dots & w_{N1} \\ \vdots & \ddots & \vdots \\ w_{1N} & \dots & w_{NN} \end{bmatrix} \begin{bmatrix} x_{11} & x_{12} & \dots & x_{1N} \\ x_{21} & x_{22} & \dots & x_{2N} \\ \vdots & \vdots & \ddots & \vdots \\ x_{M1} & x_{M2} & \dots & x_{MN} \end{bmatrix}^T$$

$$\mathbf{Y}^T = \mathbf{W} \mathbf{X}^T$$

The Fourier Transform

$$Y_k = \sum_{n=1}^N W_{kn} X_n \quad W_{kn} = e^{-j2\pi kn/N}$$

$$\mathbf{W} = \begin{bmatrix} e^{-j2\pi} & e^{-j4\pi} & \dots & e^{-j2\pi N} \\ e^{-j4\pi} & e^{-j8\pi} & \dots & e^{-j4\pi N} \\ \vdots & \vdots & \ddots & \vdots \\ e^{-j2\pi M} & e^{-j4\pi M} & \dots & e^{-j2\pi MN} \end{bmatrix}$$

(Independence between components is assumed)

Recap: Non-causal Wiener filtering

$$H(f) = \frac{S_y(f)}{S_y(f) + S_d(f)}$$

Observation $S_x(f)$
 Noise Power $S_d(f)$
 Ideal Signal $S_y(f)$

$x[n]$ - observation
 $y[n]$ - ideal signal
 $d[n]$ - noise component

Filtered signal:
 $S_{filt}(f) = S_x(f) \cdot H(f)$

BSS is a transform?

- Like Fourier, we decompose into components by transforming the observations into another vector space which maximises the separation between interesting (*signal*) and unwanted (*noise*).
- Unlike Fourier, separation is not based on frequency- It's based on *independence*
- Sources can have the same frequency content
- No assumptions about the signals (other than they are *independent* and *linearly* mixed)
- So you can filter/separate in-band noise/signals with BSS

Principal Component Analysis

- Second order decorrelation = *independence*
- Find a set of orthogonal axes in the data (independence metric = variance)
- Project data onto these axes to *decorrelate*
- *Independence* is forced onto the data through the orthogonality of axes
- Conventional noise / signal separation technique

Singular Value Decomposition

Decompose observation $X=AZ$ into....

$$X=USV^T$$

- S is a diagonal matrix of singular values with elements arranged in descending order of magnitude (the singular spectrum)
- The columns of V are the eigenvectors of $C=X^T X$ (the orthogonal subspace ... $dot(v_i, v_j)=0$) ... they 'demix' or rotate the data
- U is the matrix of projections of X onto the eigenvectors of C ... the 'source' estimates

Singular Value Decomposition

Decompose observation $X=AZ$ into....

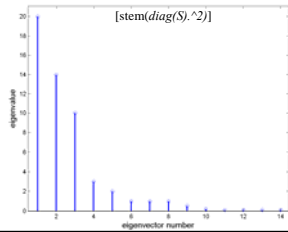
$$X=USV^T$$

$$\begin{bmatrix} x_{11} & x_{12} & \dots & x_{1N} \\ x_{21} & x_{22} & \dots & x_{2N} \\ \vdots & \vdots & \ddots & \vdots \\ \vdots & \vdots & \ddots & \vdots \\ x_{M1} & x_{M2} & \dots & x_{MN} \end{bmatrix} = \begin{bmatrix} u_{11} & u_{12} & \dots & u_{1N} \\ u_{21} & u_{22} & \dots & u_{2N} \\ \vdots & \vdots & \ddots & \vdots \\ \vdots & \vdots & \ddots & \vdots \\ u_{M1} & u_{M2} & \dots & u_{MN} \end{bmatrix} + \begin{bmatrix} s_{11} & 0 & \dots & 0 \\ 0 & s_{22} & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & s_{NN} \\ \vdots & \vdots & \ddots & \vdots \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & 0 \end{bmatrix} + \begin{bmatrix} v_{11} & \dots & v_{1N} \\ v_{21} & \dots & v_{2N} \\ \vdots & \vdots & \vdots \\ v_{N1} & \dots & v_{NN} \end{bmatrix}$$

$X = U \quad S \quad V^T$

Eigenspectrum of decomposition

- S = singular matrix ... zeros except on the leading diagonal
- S_{ij} ($i=j$) are the *eigenvalues*^{1/2}
- Placed in order of descending magnitude
- Correspond to the magnitude of projected data along each *eigenvector*
- Eigenvectors are the axes of maximal variation in the data



Eigenspectrum =
Plot of eigenvalues

Variance = power
(analogous to Fourier
components in power spectra)

SVD: Method for PCA

A routine for performing SVD is as follows:

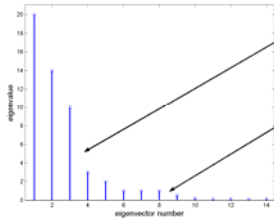
1. Find the N non-zero eigenvalues, λ_i of the matrix $C = X^T X$ and form a non-square diagonal matrix S by placing the square roots $s_i = \sqrt{\lambda_i}$ of the N eigenvalues in descending order of magnitude on the leading diagonal and setting all other elements of S to zero.
2. Find the orthogonal eigenvectors of the matrix $X^T X$ corresponding to the obtained eigenvalues, and arrange them in the same order. this ordered collection of column-vectors forms the matrix V .
3. Find the first N column-vectors of the matrix U : $u_i = s_i^{-1} X v_i$ ($i = 1 : N$). Note that s_i^{-1} are the elements of S^{-1} .
4. Add the rest of $M - N$ vectors to the matrix U using the Gram-Schmidt orthogonalization process (see appendix 15.9.2).

SVD noise/signal separation

To perform SVD filtering of a signal, use a truncated SVD decomposition (using the first p eigenvectors)

$$Y = US_p V^T$$

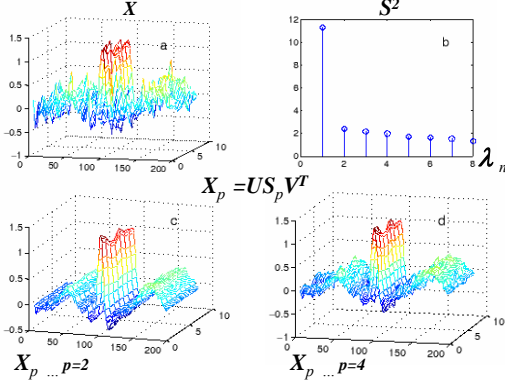
[Reduce the dimensionality of the data by discarding noise projections $S_{\text{noise}} = 0$
Then reconstruct the data with just the signal subspace]



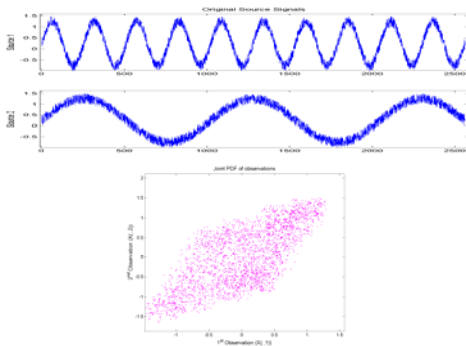
Most of the signal is contained in the first few principal components.

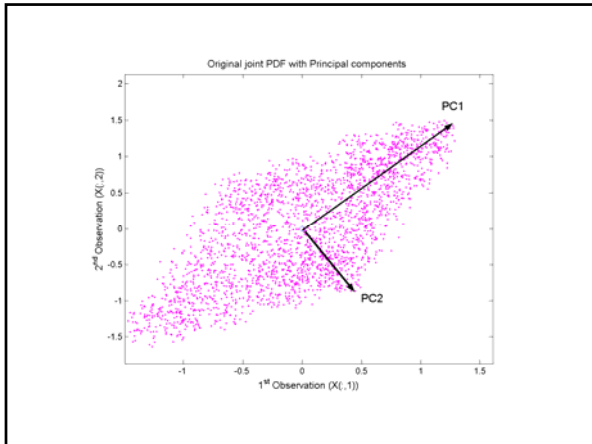
Discarding these and projecting back into the original observation space effects a noise-filtering or a noise/signal separation

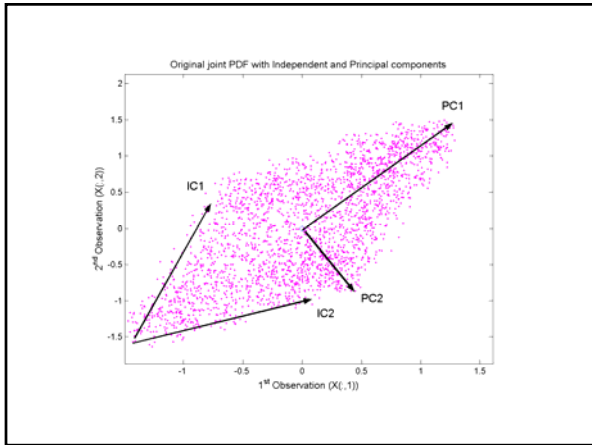
Real data



Two dimensional example







Independent Component Analysis

As in PCA, we are looking for N different vectors onto which we can project our observations to give a set of N maximally independent signals (*sources*)

output data (discovered sources) dimensionality = dimensionality of observations

Instead of using *variance* as our independence measure (i.e. decorrelating) as we do in PCA, we use a measure of how statistically independent the sources are.

ICA: The basic idea ...

Assume underlying source signals (Z) are independent.

Assume a linear mixing matrix (A)... $X^T = AZ^T$

in order to find $Y (\approx Z)$, find W , ($\approx A^{-1}$) ...

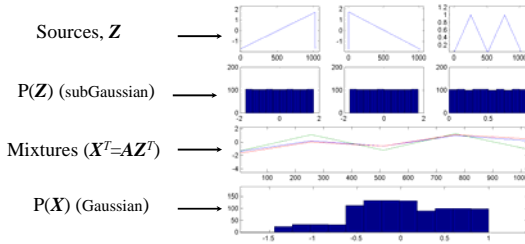
$$Y^T = WX^T$$

How? Initialise W & iteratively update W to minimise or maximise a cost function that measures the (statistical) *independence* between the columns of the Y^T .

Non-Gaussianity \Rightarrow statistical independence?

From the *Central Limit Theorem*,

- add enough independent signals together, \rightarrow Gaussian PDF



Recap: Moments of a distribution

$$\mu_x = E\{x\} = \int_{-\infty}^{+\infty} x p_x(x) dx \quad \left| \quad \hat{\mu}_x = \frac{1}{M} \sum_{i=1}^M x_i$$

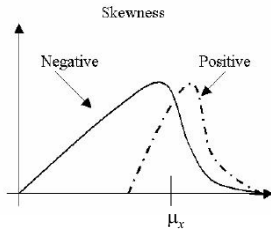
$$\sigma_x^2 = E\{(x - \mu_x)^2\} = \int_{-\infty}^{+\infty} (x - \mu_x)^2 p_x(x) dx \quad \hat{\sigma}^2(x) = \frac{1}{M} \sum_{i=1}^M (x_i - \hat{\mu}_x)^2$$

$$\sigma(x) = \sqrt{\sigma^2}$$

$$v_n = E\{(x - \mu_x)^n\} = \int_{-\infty}^{+\infty} (x - \mu_x)^n p_x(x) dx$$

Higher order moments (3rd-skewness)

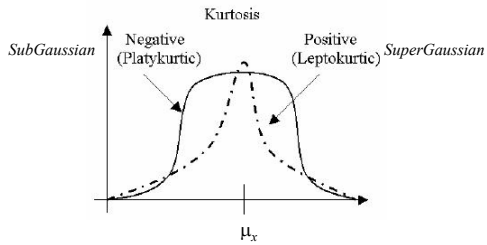
$$\hat{\zeta}(\mathbf{x}) = \frac{1}{M} \sum_{i=1}^M \left[\frac{x_i - \hat{\mu}_x}{\hat{\sigma}} \right]^3$$



Higher order moments (4th-kurtosis)

$$\hat{\kappa}(\mathbf{x}) = \frac{1}{M} \sum_{i=1}^M \left[\frac{x_i - \hat{\mu}_x}{\hat{\sigma}} \right]^4$$

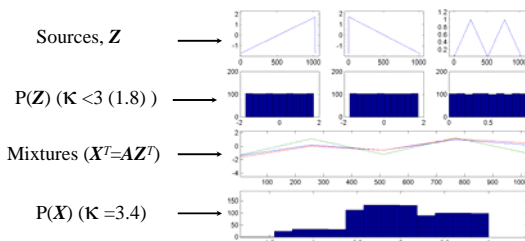
Gaussians are *mesokurtic* with $\kappa = 3$



Non-Gaussianity \Rightarrow statistical independence?

Central Limit Theorem: add enough independent signals together,
 \rightarrow Gaussian PDF ($\kappa = 3$)

\therefore make data components non-Gaussian to find independent sources



Recall – trying to estimate \mathbf{W}

Assume underlying source signals (\mathbf{Z}) are independent.

Assume a linear mixing matrix (\mathbf{A})... $\mathbf{X}^T = \mathbf{AZ}^T$

in order to find \mathbf{Y} ($\approx \mathbf{Z}$), find \mathbf{W} , ($\approx \mathbf{A}^{-T}$) ...

$$\mathbf{Y}^T = \mathbf{W}\mathbf{X}^T$$

Initialise \mathbf{W} & iteratively update \mathbf{W} with gradient descent to maximise kurtosis.

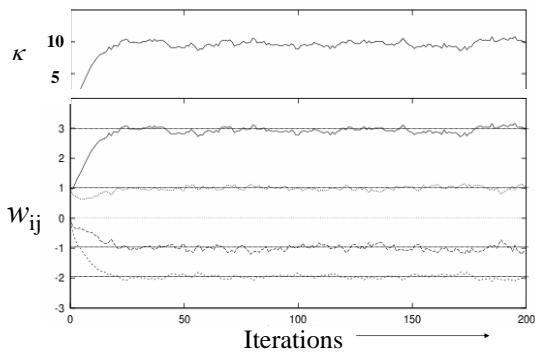
Gradient descent to find \mathbf{W}

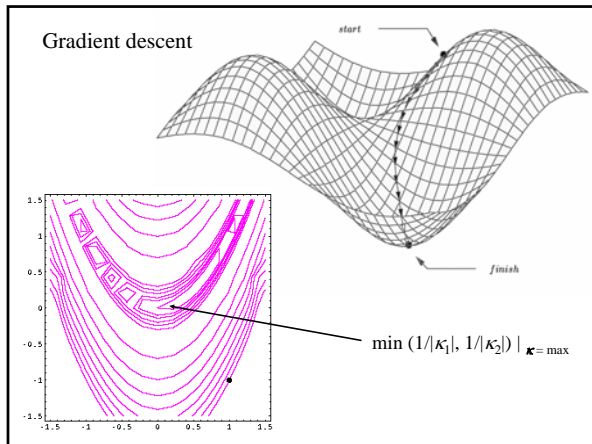
- Given a cost function, ξ , we update each element of \mathbf{W} (w_{ij}) at each step, τ ,

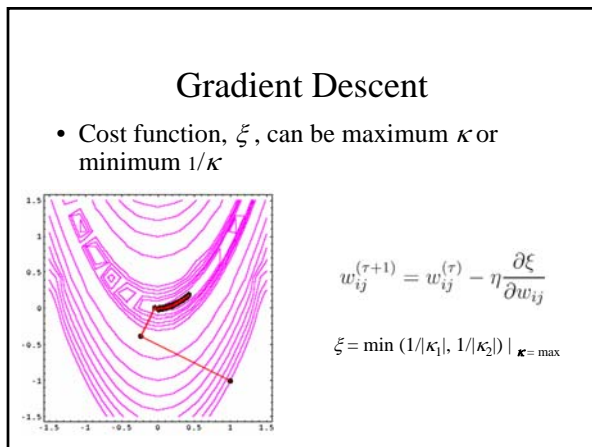
$$w_{ij}^{(\tau+1)} = w_{ij}^{(\tau)} - \eta \frac{\partial \xi}{\partial w_{ij}}$$

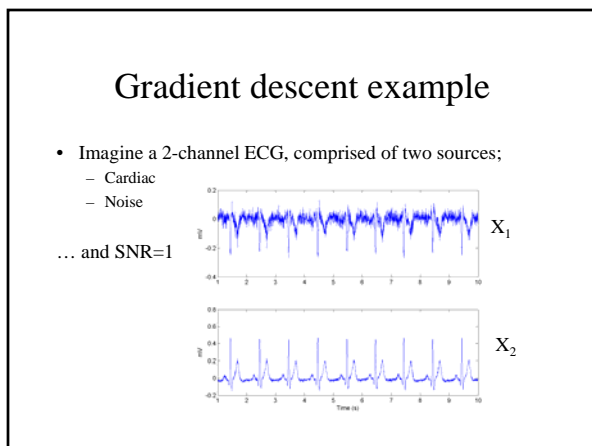
- ... and recalculate cost function
- (η is the learning rate (~ 0.1), and speeds up convergence.)

Weight updates to find: $\mathbf{W} = \begin{bmatrix} 1 & 3 \\ -2 & -1 \end{bmatrix}$
(Gradient ascent)

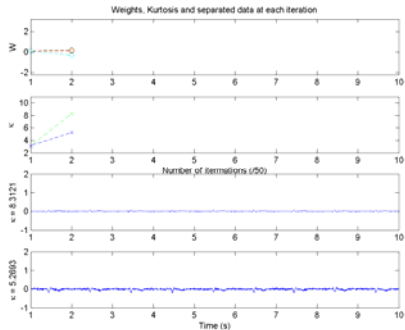








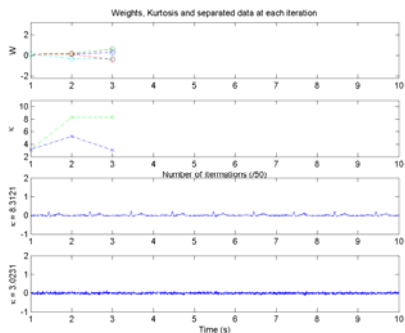
Iteratively update W and measure κ



Y_1

Y_2

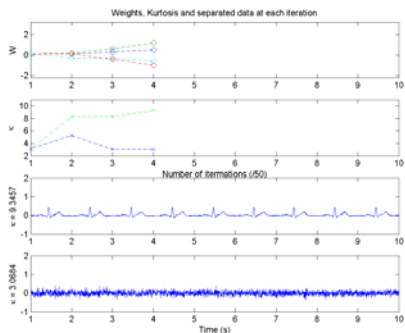
Iteratively update W and measure κ



Y_1

Y_2

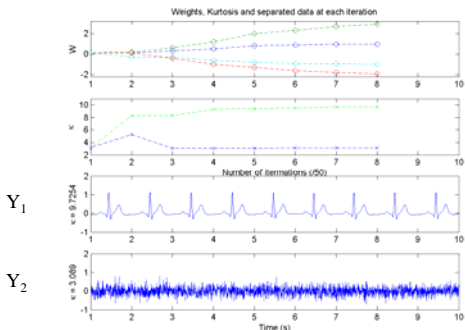
Iteratively update W and measure κ



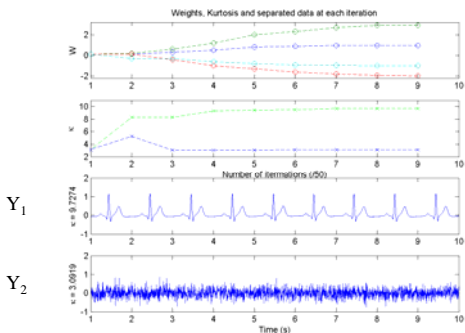
Y_1

Y_2

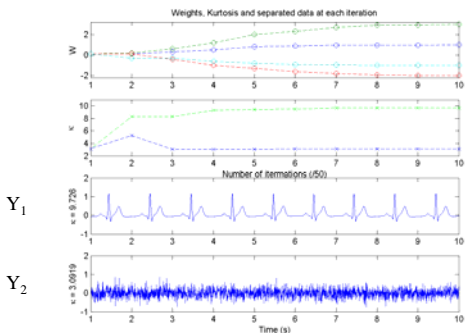
Iteratively update W and measure κ



Iteratively update W and measure κ



Maximized κ for non-Gaussian signal



Outlier insensitive ICA cost functions

Measures of statistical independence

In general we require a measure of statistical independence which we maximise between each of the N components.

Non-Gaussianity is one approximation, but sensitive to small changes in the distribution tail.

Other measures include:

- Mutual Information I ,
- Entropy (Negentropy, \mathcal{J})... and
- Maximum (Log) Likelihood $\mathcal{L}(\mathbf{W})$

(Note: all are related to κ)

Entropy-based cost function

Kurtosis is highly sensitive to small changes in distribution tails.

A more robust measures of Gaussianity is based on differential entropy $H(\mathbf{y})$,

$$H(\mathbf{y}) = - \int P(\mathbf{y}) \log_2 P(\mathbf{y}) d\mathbf{y}$$

... *negentropy*:

$$\mathcal{J}(\mathbf{y}) = H(\mathbf{y}_{\text{gauss}}) - H(\mathbf{y})$$

where $\mathbf{y}_{\text{gauss}}$ is a Gaussian variable with the same covariance matrix as \mathbf{y} . $\mathcal{J}(\mathbf{y})$ can be estimated from kurtosis ...

$$\mathcal{J}(\mathbf{y}) \approx \frac{1}{12} E\{y^3\}^2 + \frac{1}{48} \kappa(y)^2$$

Entropy: measure of randomness- Gaussians are maximally random

Minimising Mutual Information

Mutual information (MI) between two vectors \mathbf{x} and \mathbf{y} :

$$I = H_x + H_y - H_{xy}$$

always non-negative and zero if variables are independent ...
therefore we want to minimise MI.

MI can be re-written in terms of negentropy ...

$$I(y_1, y_2, \dots, y_m) = c - \sum_{i=1}^m J(y_i)$$

where c is a constant.

... differs from negentropy by a constant and a sign change

Independent source discovery using Maximum Likelihood

Generative latent variable modelling N observables, \mathbf{X} ...
from N sources, z_i through a linear mapping $\mathbf{W}=w_{ij}$

Latent variables assumed to be independently distributed

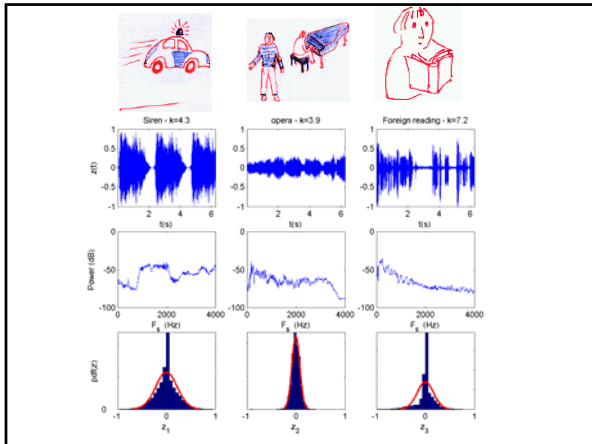
Find elements of \mathbf{W} by gradient ascent $\Delta w_{ij} = \eta \frac{\partial \mathbf{L}}{\partial w_{ij}}$
- iterative update by

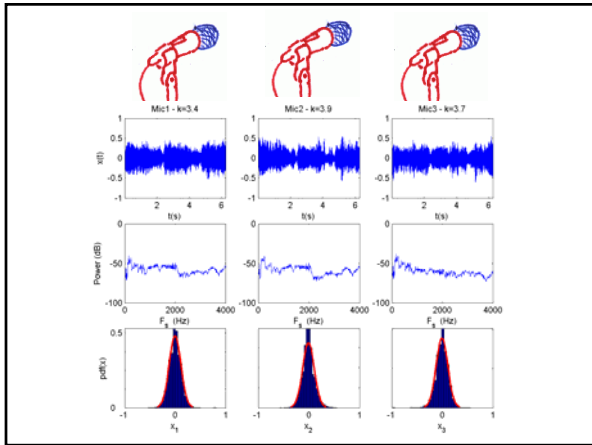
where η is some learning rate (const) ... and
 $\mathbf{L}(\mathbf{W})$ is our objective *cost function*, the **log likelihood**

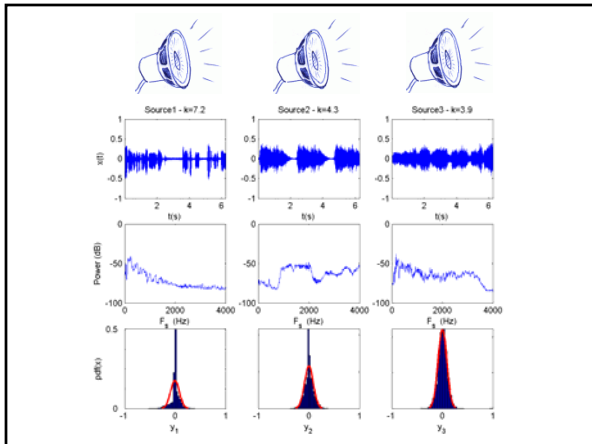
$$\log_2 P(\mathbf{x}^m | \mathbf{A}) = \log_2 \det \mathbf{A} + \sum_i \log_2 p_i(a_{ij} \mathbf{x}_j)$$

The cocktail party problem revisited

... some real examples using ICA







Observations

Separation of mixed observations into source estimates is excellent ... apart from:

- Order of sources has changed
- Signals have been scaled

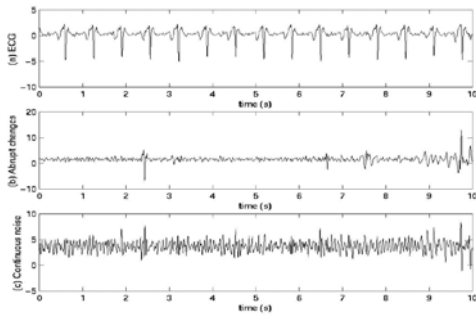
Why? ... In $X^T = AZ^T$, insert a permutation matrix B ...

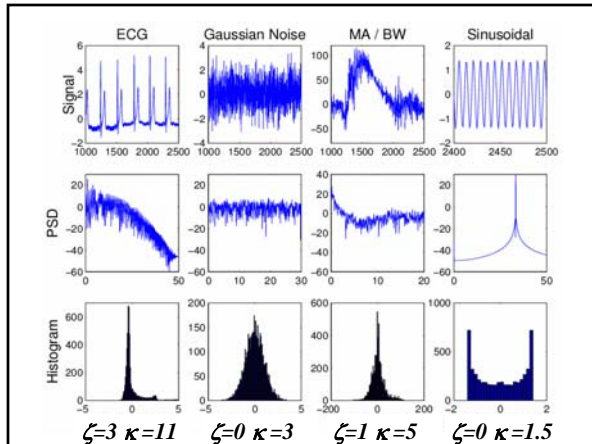
$$X^T = AB B^{-1} Z^T \Rightarrow B^{-1} Z^T \dots = \text{sources with different col. order.}$$

\Rightarrow sources change by a scaling $A \mapsto AB$

... ICA solutions are order and scale independent because κ is dimensionless

Separation of sources in the ECG

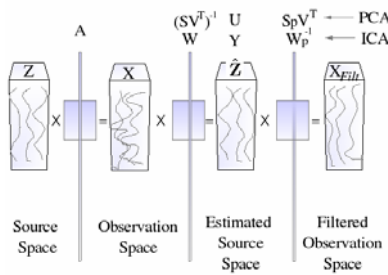




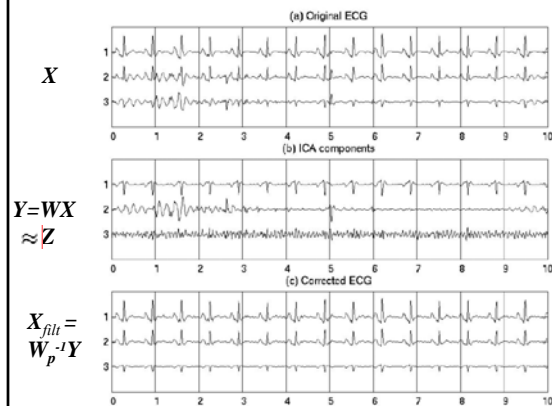
Transformation inversion for filtering

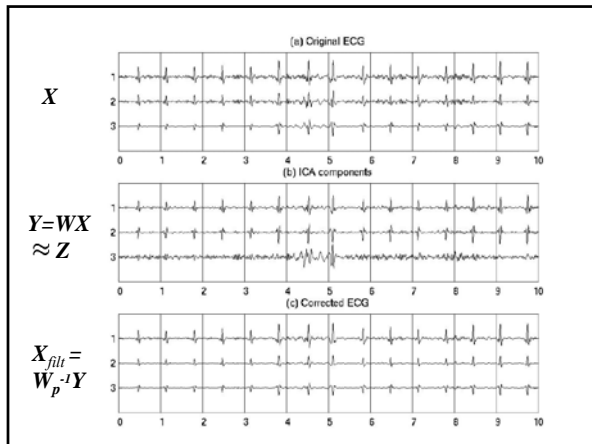
- Problem - can never know if sources are really reflective of the actual source generators - no gold standard
- De-mixing might alter the clinical relevance of the ECG features
- Solution: Identify unwanted sources, set corresponding (p) columns in W^{-1} to zero (W_p^{-1}), then multiply back through to remove 'noise' sources and transform back into original observation space.

Transformation inversion for filtering



Real data





Summary

- PCA is good for **Gaussian** noise separation
- ICA is good for **non-Gaussian** 'noise' separation
- PCs have obvious meaning - highest *energy* components
- ICA - derived sources : **arbitrary scaling/inversion & ordering**
 need energy-independent heuristic to identify signals / noise
- Order of ICs change - IC space is **derived** from the data.
 - PC space only changes if SNR changes.
- ICA assumes **linear mixing matrix**
- ICA assumes **stationary** mixing
- De-mixing performance is function of lead position
- ICA requires as many sensors (ECG leads) as *sources*
- Filtering - discard certain dimensions then *invert transformation*
- In-band noise can be removed - unlike Fourier!

