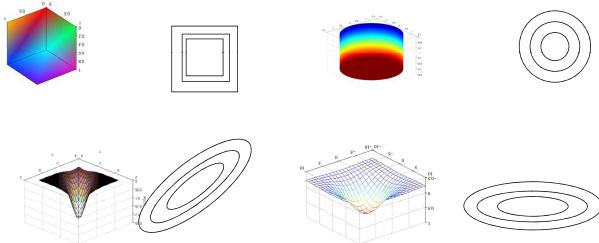


- Which of the above represent uncorrelated RVs?
- Which of the above represent independent RVs?



## Independence

$$\boxed{E[f(x)g(y)] = E[f(x)]E[g(y)] \text{ for all } f(), g()}$$

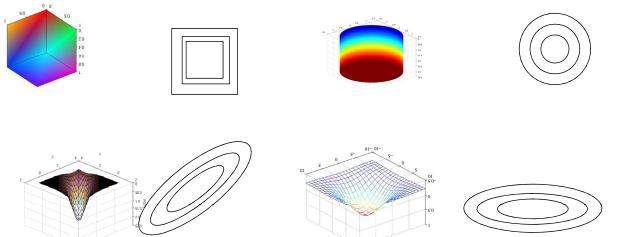
- The average value of any function  $X$  is the same regardless of the value of  $Y$
- Independence: Two random variables  $X$  and  $Y$  are independent iff:

## A brief review of basic probability

- $E[X|Y] = E[X]$
- The average value of  $X$  is the same regardless of the value of  $Y$
- $\leftarrow$  The average value of  $X$  is the same regardless of the joint probability of  $(X,Y)$
- $P(X,Y) = P(X)P(Y)$
- Individual probabilities
- Their joint probability equals the product of their individual probabilities
- Independence: Two random variables  $X$  and  $Y$  are independent iff:

## A brief review of basic probability

- Which of the above represent uncorrelated RVs?



## Uncorrelatedness

- The average value of  $X$  is the same regardless of the value of  $Y$
- $E[XY] = E[X]E[Y]$
- I.e. one instance of  $(X,Y)$
- Instance of  $Y$
- Setup: Each draw produces one instance of  $X$  and one product of their individual averages
- The average value of the product of the variables equals the product of their individual averages
- Uncorrelated iff:
- Uncorrelated: Two random variables  $X$  and  $Y$  are

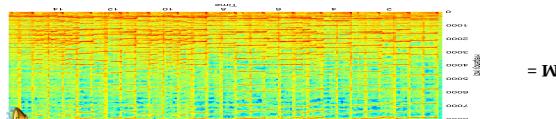
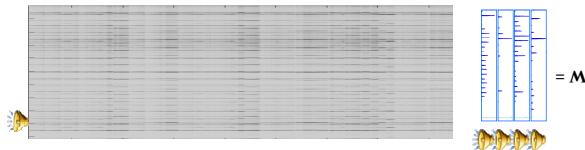
## A brief review of basic probability

Class 20, 8 Nov 2012

Instructor: Bhiksha Raj

# Independent Component Analysis

- Projected Spectrogram =  $P * M$
- $P = W(W^T W)^{-1} W^T$



Projection: multiple notes

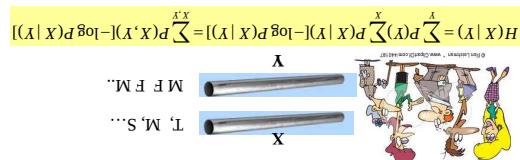
Oneward.

- Entropies of X and Y if they are independent
- Joint entropy of X and Y is the sum of the entropies of X and Y ( $H(X, Y) = H(X) + H(Y)$ )
- Conditional entropy of X given Y is the average of the number of bits to transmit a symbol Y given X ( $H(X|Y) = \sum_x p(x) \sum_y p(y|x) (-\log p(y|x))$ )
- Conditional entropy of X =  $H(X)$  if X is independent of Y

A brief review of basic info. theory

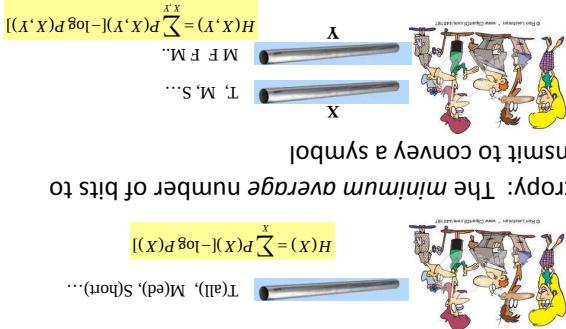
■ Averaged over all values of X and Y

- After symbol Y has already been conveyed, the number of bits to transmit to convey a symbol X, given Y, is  $H(X|Y) = \sum_x p(x|y) \sum_y p(y) (-\log p(x|y))$
- Conditional Entropy: The minimum average



A brief review of basic info. theory

- Entropy: The minimum average number of bits to convey a symbol to convey a symbol
- Joint entropy: The minimum average number of bits to convey sets (pairs here) of symbols

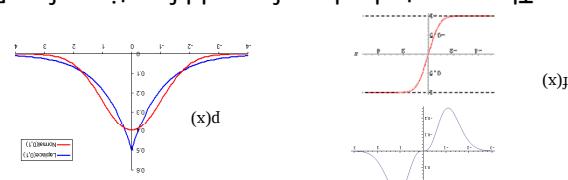


A brief review of basic probability

■  $E[f(x)] = 0$  if  $f(x)$  is odd symmetric

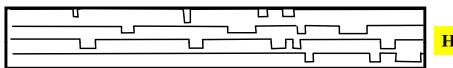
- The PDF is of the RV is symmetric around 0
- The RV is 0 mean

■ The expected value of an odd function of an RV is 0



A brief review of basic probability

- The rows of  $H$  are uncorrelated
- $h_i^T h_j = 0$  for all  $i \neq j$
- When one note occurs, the other does not
- For our problem, lets consider the "truth".



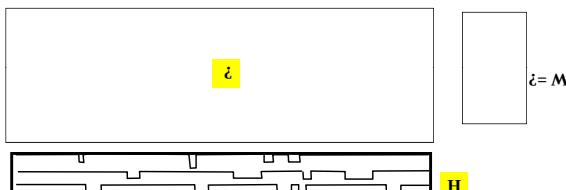
- Unconstrained
- For any  $W, H$  that minimizes the error,  $W^* = WA$ ,  $H^* = A^{-1}H$
- also minimizes the error for any invertible  $A$

$$W, H = \arg \min_{W, H} \| M - WH \|^F$$

A least squares solution

### transcription

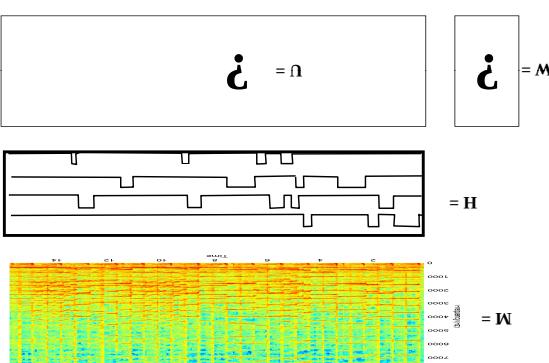
- Must ideally find the notes corresponding to the transcription
- Given  $H$ , estimate  $W$  to minimize error
- $M \sim WH$  is an approximation



Giving the other way..

$$W = \text{Pinv}(V) M$$

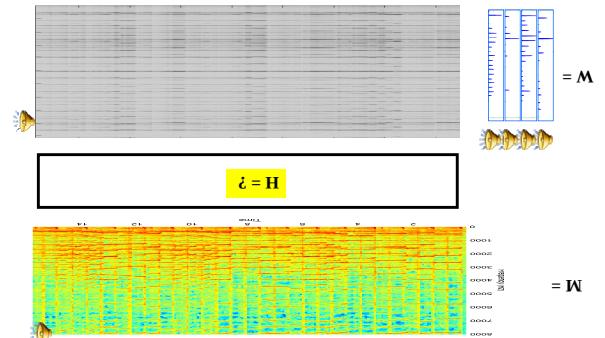
$$M \sim WH$$



How about the other way?

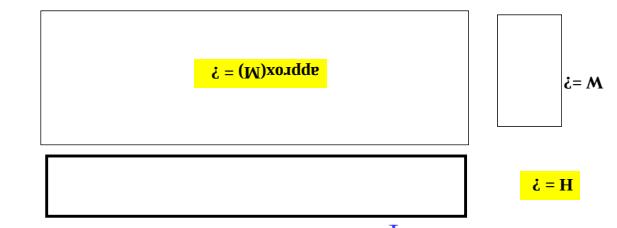
$$H = \text{Pinv}(W) M$$

$$M \sim WH$$



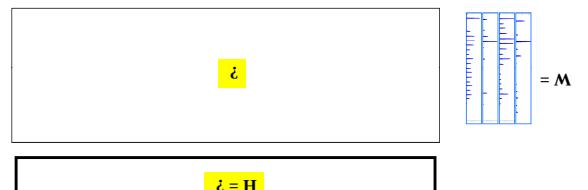
We're actually computing a score

- Must learn both the notes and their approximate  $M$
- Ideally, must learn both the notes and their transcription!



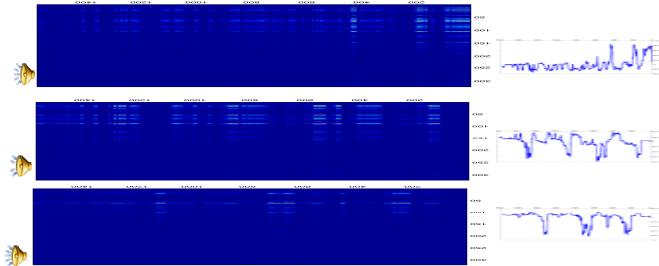
When both parameters are unknown

- Must ideally find transcription of given notes
- Given  $W$ , estimate  $H$  to minimize error
- $M \sim WH$  is an approximation



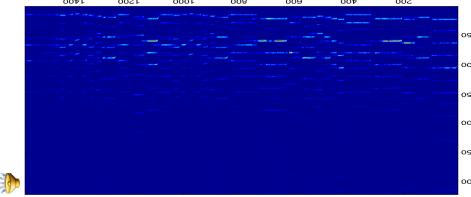
So what are we doing here?

- The first three "notes" and their contributions
- The spectrograms of the notes are statistically uncorrelated



So how does that work?

- There are 12 notes in the segment, hence we try to estimate 12 notes..



So how does that work?

- Minimize least squares error with the constraint that the rows of  $\mathbf{H}$  are length 1 and orthogonal to one another
- $\mathbf{W}, \mathbf{H} = \arg \min_{\mathbf{W}, \mathbf{H}} \| \mathbf{M} - \mathbf{WH} \|_F^2 + \sum_i \| \mathbf{h}_i \|_2^2 + \sum_{i \neq j} \mathbf{h}_i^T \mathbf{h}_j$
- is identical to

$$\mathbf{H} = \arg \max_{\mathbf{H}} \text{trace}(\text{Correlation}(\mathbf{M}^T \mathbf{H}^T \mathbf{H}) - \text{trace}(\mathbf{AH}^T \mathbf{H}))$$

## Equivalences

- Simply requiring the rows of  $\mathbf{H}$  to be orthonormal gives us that  $\mathbf{H}$  is the set of eigenvectors of the data in  $\mathbf{M}^T$

$$\text{Correlation}(\mathbf{M}^T) \mathbf{H} = \mathbf{H}\mathbf{A}$$

- Differentiating and equating to 0

$$\mathbf{H} = \arg \max_{\mathbf{H}} \text{trace}(\text{Correlation}(\mathbf{M}^T \mathbf{H}^T \mathbf{H}) - \text{trace}(\mathbf{AH}^T \mathbf{H}))$$

- Constraint: every row of  $\mathbf{H}$  has length 1

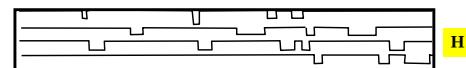
## Finding the notes

- $\mathbf{H} = \arg \max_{\mathbf{H}} \text{trace}(\text{Correlation}(\mathbf{M}^T) \mathbf{H}^T \mathbf{H})$
- $\mathbf{H} = \arg \min_{\mathbf{H}} \text{trace}(\text{Correlation}(\mathbf{M}^T)(\mathbf{I} - \mathbf{H}^T \mathbf{H}))$
- $\mathbf{H} = \arg \min_{\mathbf{H}} \text{trace}(\mathbf{M}^T \mathbf{M}(\mathbf{I} - \mathbf{H}^T \mathbf{H}))$
- $\mathbf{H} = \arg \min_{\mathbf{H}} \text{trace}(\mathbf{M}^T (\mathbf{I} - \mathbf{H}^T \mathbf{H}) \mathbf{M}^T)$
- Could also be rewritten as
  - Only  $\mathbf{HH}^T = \mathbf{I}$
- Note  $\mathbf{H}^T \mathbf{H} = \mathbf{I}$

$$\mathbf{H} = \arg \min_{\mathbf{H}} \| \mathbf{M} - \mathbf{MH}^T \mathbf{H} \|_F^2$$

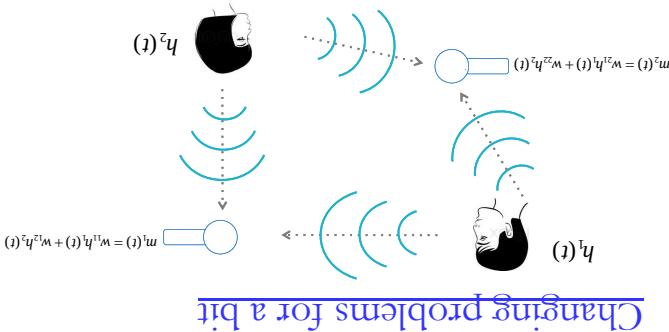
Finding the notes

- $\mathbf{H} = \arg \min_{\mathbf{H}} \| \mathbf{M} - \mathbf{WH} \|_F^2$  Constraint:  $\text{Rank}(\mathbf{H}) = 4$
- $\mathbf{W} = \mathbf{M} \mathbf{H}^T \mathbf{H}$
- $\mathbf{W} = \mathbf{M} \text{pinv}(\mathbf{H}) = \mathbf{M} \mathbf{H}^T$
- Projecting  $\mathbf{M}$  onto  $\mathbf{H}$
- $\text{pinv}(\mathbf{H}) = \mathbf{H}^T$
- Normalizing all rows of  $\mathbf{H}$  to length 1
- Assume:  $\mathbf{HH}^T = \mathbf{I}$

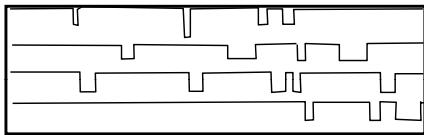


A least squares solution

- Each recorded signal is a mixture of both signals
- Recorded by two microphones
- Two people speak simultaneously

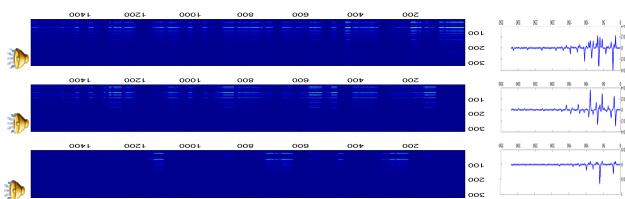


- Not strictly true, but still ..
- Playing independently of one another
- Or, in a multi-instrument piece, instruments are dependent on what else is playing
- Assume: The "transcription" of one note does not



What else can we look for?

- There are 12 notes in the segment, hence we try to estimate 12 notes..



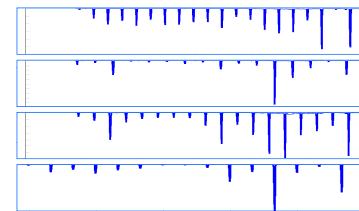
So how does that work?

- Impose statistical independence constraints on decomposition

$$\mathbf{W}, \mathbf{H} = \arg \min_{\mathbf{W}, \mathbf{H}} \| \mathbf{M} - \mathbf{WH} \|_F^2 + \lambda (\text{rows of } \mathbf{H} \text{ are independent})$$

Formulating it with Independence

- More generally, simple orthogonality will not give us the desired solution
- Harmonica continues to resonate to previous note
- Note occurs concurrently
- Overlapping frequencies



Our notes are not orthogonal

- Solving the above, with the constraints that the columns of  $\mathbf{W}$  are orthonormal gives you the eigen vectors of the data in  $\mathbf{M}$
- $\mathbf{W} = \arg \max_{\mathbf{W}} \text{trace}(\mathbf{W}^T \mathbf{W} \text{Correlation}(\mathbf{M})) - \text{trace}(\mathbf{AW}^T \mathbf{W})$
- $\text{Correlation}(\mathbf{M}) \mathbf{W} = \mathbf{AW}$

$$\mathbf{W} = \arg \min_{\mathbf{W}} \| \mathbf{M} - \mathbf{W}^T \mathbf{WM} \|_F^2$$

- Can find  $\mathbf{W}$  instead of  $\mathbf{H}$

Finding the notes

- In reality, we only want this to be a diagonal matrix, but we'll make it identity
- $\mathbf{X} = \mathbf{CM}$
- Estimate a  $\mathbf{C}$  such that  $\mathbf{CM}$  is uncorrelated
- Independence  $\rightarrow$  Uncorrelatedness

$$\mathbf{H} = \mathbf{AM} \quad \mathbf{A} = \mathbf{BC}$$

$$\mathbf{H} = \mathbf{BCM}$$

## Emulating Independence

- $\mathbf{m}_i$  are the columns of  $\mathbf{M}$

$$\mathbf{m}_i = \mathbf{m}_i - \bar{\mathbf{m}}$$

$$\bar{\mathbf{m}}_i = \frac{\text{cols}(\mathbf{M})}{i} \sum_{j=1}^i \mathbf{m}_j$$

- First step of ICA: Set the mean of  $\mathbf{M}$  to 0
- $\mathbf{E}[\mathbf{H}] = \mathbf{AE}[\mathbf{M}] = \mathbf{AO} = \mathbf{0}$
- If  $\text{mean}(\mathbf{M}) = 0 \Rightarrow \text{mean}(\mathbf{H}) = 0$

$$\mathbf{M} = \mathbf{WH} \quad \mathbf{H} = \mathbf{AM}$$

- Usual to assume zero mean processes
- Otherwise, some of the math doesn't work well

## Zero Mean

- The fourth order moments are independent
- $\mathbf{h}_i$  and  $\mathbf{h}_j$  are the  $i$ th and  $j$ th components of any vector in  $\mathbf{H}$
- $\mathbf{E}[\mathbf{h}_i \mathbf{h}_j] = \mathbf{E}[\mathbf{h}_i] \mathbf{E}[\mathbf{h}_j]$
- The rows of  $\mathbf{H}$  are uncorrelated

$$\mathbf{H} = \begin{bmatrix} \text{Signal from speaker 1} \\ \text{Signal at mic 1} \\ \text{Signal from speaker 2} \\ \text{Signal at mic 2} \end{bmatrix}$$

## Emulating Independence

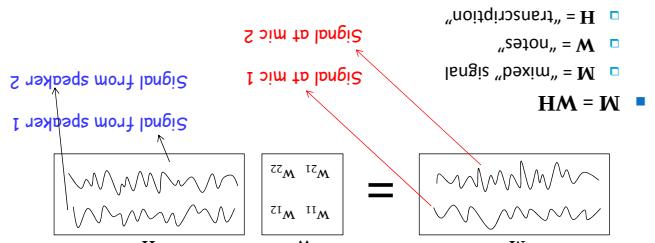
- Estimate  $\mathbf{A}$  such that the components of  $\mathbf{AM}$  are statistically independent
- A is the unmixing matrix
- Ensure that the components of the vectors in the estimated  $\mathbf{H}$  are statistically independent
- Given only  $\mathbf{M}$  estimate  $\mathbf{H}$
- $\mathbf{H} = \mathbf{W}^{-1}\mathbf{M} = \mathbf{AM}$

$$\mathbf{M} = \mathbf{WH} \quad \mathbf{W} = \begin{bmatrix} w_{11} & w_{12} \\ w_{21} & w_{22} \end{bmatrix} \quad \mathbf{H} = \begin{bmatrix} \text{Signal from speaker 1} \\ \text{Signal at mic 1} \\ \text{Signal from speaker 2} \\ \text{Signal at mic 2} \end{bmatrix}$$

## Imposeing Statistical Constraints

- Given only  $\mathbf{M}$  estimate  $\mathbf{H}$
- $\mathbf{H} = \mathbf{W}^{-1}\mathbf{M} = \mathbf{AM}$
- Ensure that the components of the vectors in the estimated  $\mathbf{H}$  are statistically independent
- A is the unmixing matrix
- Multiple approaches.

- Given only  $\mathbf{M}$  estimate  $\mathbf{H}$
- $\mathbf{H} = \mathbf{W}^{-1}\mathbf{M} = \mathbf{AM}$
- Ensure that the components of the vectors in the estimated  $\mathbf{H}$  are statistically independent
- A is the unmixing matrix
- Multiple approaches.



## Imposeing Statistical Constraints

- Create a matrix of fourth moment terms that would be diagonal were the rows of  $\mathbf{H}$  independent and diagonalize it
  - A good candidate
  - Good because it incorporates the energy in all rows of  $\mathbf{H}$
  - The fourth moments of  $\mathbf{H}$  have the form:
  - If the rows of  $\mathbf{H}$  were independent
  - $E[h_1 h_2 h_3] = E[h_1] E[h_2] E[h_3]$
  - $D = E[h_1^2 h_2^2 h_3^2]$
  - i.e.
  - Where
  - $d_{ij} = E[h_i^2 h_j^2]$
  - $D = \begin{bmatrix} d_{11} & d_{12} & d_{13} & .. \\ d_{21} & d_{22} & d_{23} & .. \\ d_{31} & d_{32} & d_{33} & .. \\ .. & .. & .. & .. \end{bmatrix}$
  - Solution: Compute  $\mathbf{B}$  such that the fourth moments of  $\mathbf{H} = \mathbf{BX}$
  - While ensuring that  $\mathbf{B}$  is Unitary
  - $\mathbf{B}$  are the columns of  $\mathbf{H}$
  - Assuming  $\mathbf{h}$  is real, else replace transpose with Hermitian

- Uncorrelated  $\Leftrightarrow$  Independent
  - Whitening merely ensures that the resulting signals are uncorrelated, i.e.,
  - $E[x_i x_j] = 0 \text{ if } i \neq j$
  - This does not ensure higher order moments are also decoupled, e.g. it does not ensure that decoupled, e.g. it does not ensure that components?
  - Will multiplying  $X$  by  $B$  re-correlate the components?
  - Not if  $B$  is unitary
  - Since the rows of  $H$  are uncorrelated
  - Because they are fourth order moments
  - Let's explicitly decouple the four independent RVs
  - This is one of the signatures of independent RVs
  - SO we want to find a unitary matrix
  - $BB^T = B^T B = I$
  - $XX^T = I$
  - $X = CM$
  - $H$
  - $H^T$
  - Diagonal + rank-1 matrix
  - $H = B C$
  - $A = B C$
  - $H = B C$

## Decorating

Uncorrected ≡ Independent

- $H = A + B + C$

## Decorrelation



- Minimize the above to obtain  $\mathbf{B}$

$$J(\mathbf{H}) = \sum_i^l H(\underline{\mathbf{h}}_i) - \log |\mathbf{W}|$$

- Ignoring  $H(\mathbf{x})$  (Const)

$$I(\mathbf{H}) = \sum_i^l H(\underline{\mathbf{h}}_i) - H(\mathbf{x}) - \log |\mathbf{B}|$$

$$\underline{I}(\mathbf{H}) = \sum_i^l H(\underline{\mathbf{h}}_i) - H(\underline{\mathbf{H}})$$

## The contrast function

$$H(\mathbf{h}) = H(\mathbf{x}) + \log |\mathbf{B}|$$

$$H(\mathbf{x}) = \int P(\mathbf{x}) \log P(\mathbf{x}) d\mathbf{x}$$

$$P(\mathbf{h}) = P_{\mathbf{x}}(\mathbf{B}^{-1}\mathbf{h}) |\mathbf{B}|^{-1}$$

- $\mathbf{x}$  is mixed signal,  $\mathbf{B}$  is the unmixing matrix
- Individual columns of the  $\mathbf{H}$  and  $\mathbf{X}$  matrices
- $\mathbf{h} = \mathbf{Bx}$

## Linear Functions

- $\mathbf{X}$  is "whitened"  $\mathbf{M}$
- With constraint:  $\mathbf{H} = \mathbf{BX}$

$$I(\mathbf{H}) = \sum_i^l H(\underline{\mathbf{h}}_i) - H(\underline{\mathbf{H}})$$

- An explicit contrast function

are independent

- Contrast function: A non-linear function that has a minimum value when the output components are independent

## The contrast function

- Can use first  $k$  columns of  $\mathbf{E}$  only if only  $k$  independent sources are expected
- In microphone array setup – only  $K < M$  sources

$$\begin{aligned} C &= S^{-1/2} E_T \\ &\text{Eigen decomposition } \mathbf{M} \mathbf{M}^T = \mathbf{E} \mathbf{S} \mathbf{E}^T \\ &E[\mathbf{x}_i \mathbf{x}_j] = E[\mathbf{x}_i] E[\mathbf{x}_j] = \delta_{ij} \text{ for centered signals} \\ &\mathbf{X} = \mathbf{CM} \end{aligned}$$

- Eliminate second-order dependence
- Normalize variance along all directions
- The mixed signal is usually "prewhitened"

## A note on pre-whitening

- Contrast functions are often only approximations too..
- F(AM)

- Define and minimize a contrast function

independent

- Contrast function: A non-linear function that has a minimum value when the output components are independent

$$\mathbf{H} = \mathbf{AM}$$

independent

- Specifically ensure that the components of  $\mathbf{H}$  are

## Ensuring Independence

- JADE: Joint Approximate Diagonalization of Eigenmatrices, J.F. Cardoso
- Diagonalizes several fourth-order moment matrices
- More effective than the procedure shown, but more computationally expensive
- Jointly diagonalizes several fourth-order moment matrices

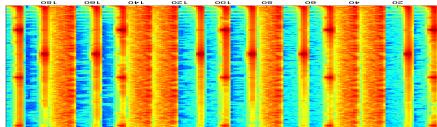
- Diagonalized to diagonalize every other fourth-order moment matrix chosen in not guaranteed to diagonalize every other fourth-order moment matrix
- There are many other ways of constructing fourth-order moment matrices that would ideally be diagonal
- Only a subset of fourth order moments are considered
- The procedure just outlined, while fully functional, has shortcomings
- The procedure just outlined, while fully functional, has

## matrices

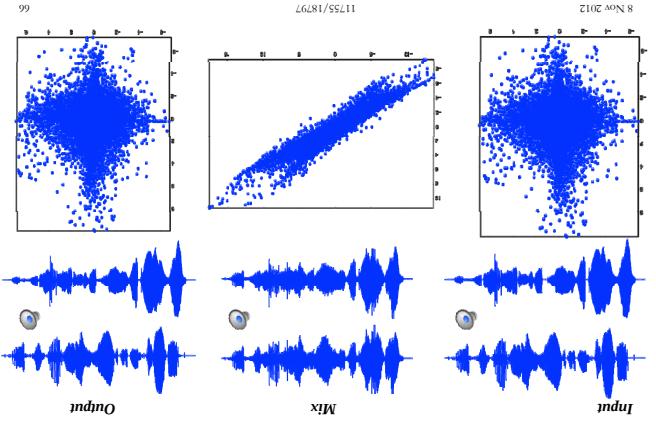
## ICA by diagonalizing moment



- Three instruments..

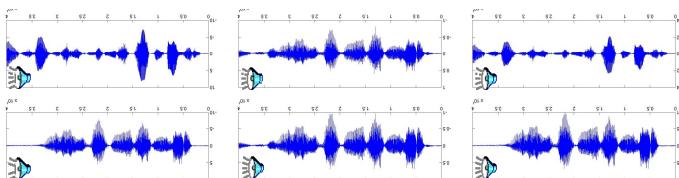


## Another Example



## Another example

- Works very well
- Natural gradient update
- speakers
- Example with instantaneous mixture of two



## So how does it work?

$$g(x) = \begin{cases} x - \tanh(x) & x \text{ is sub Gaussian} \\ x + \tanh(x) & x \text{ is super Gaussian} \end{cases}$$

- Multiple functions proposed
- Must be odd symmetric functions

## What are G and H?

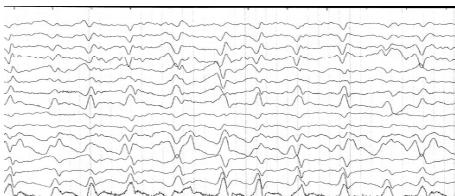
- $\Delta B = (I - g(H)f(H)^T)W$
- Cicikli-Unbehauen
- $\Delta B = (I - g(H)H^T)W$
- Natural gradient --  $f()$  = identity function
- $B = B + \eta \Delta B$
- $H = BX$
- for  $g()$  and  $f()$
- Multiple solutions under different assumptions

## Update Rules

- $\Delta B = ((B^{-1} - g(H))X)$
- Bell Sejnowski
- neural network
- $\Delta B_{ij} = f(h_i)g(h_j)$ ; -- actually assumed a recursive
- Juttien Hereraut : Online update
- $B = B + \eta \Delta B$
- $H = BX$
- for  $g()$  and  $f()$
- Multiple solutions under different assumptions

## Update Rules

- Very commonly used to enhance EEG signals
  - EEG signals are frequently corrupted by heartbeats and biohythms signals
  - ICA can be used to separate them out



## ICA for Signal Enhancement

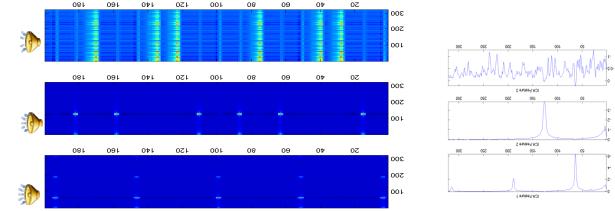
- Audio preprocessing example
  - Take a lot of audio snippets and concatenate them in a big matrix, do component analysis
  - PCA results in the DCT bases
  - ICA returns time/frequency localized sinusoids which is a better way to analyze sounds
  - DITTO for images
  - ICA returns edge filters

## Finding useful transforms with ICA

- The “bases” in PCA
  - represents the “building blocks”
  - ideally notes
  - Very successfully used
  - So can ICA be used to do the same?

## IICA for data exploration

- ## ■ Three Instruments..



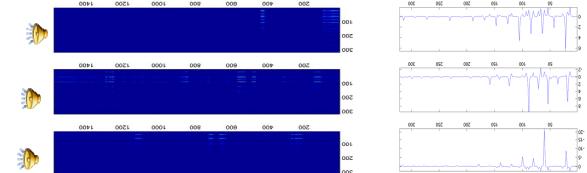
The Notes

- If one note plays, other notes are *not* playing
  - Only one note plays at a time
  - Notes are *not independent*
  - Still this didn't affect the three instruments case..
  - Not symmetric – negative values never happen
  - Note energy here
  - Assume distribution of signals is symmetric around mean
  - Continue in next class..
- What else went wrong?**

- In worse case, output are not desired signals at all..
- In the best case
- permuted order
- Outputs are scaled versions of desired signals in
- Scaling the signal does not affect independence
- Does not have sense of scaling
- Permutation invariance
- So the sources can come in any order
- Get k independent directions, but does not have a notion of the "best" direction
- Unlike PCA
- No sense of order
- Does not have sense of scaling
- Permutation invariance
- So the sources can come in any order
- Get k independent directions, but does not have a notion of the "best" direction
- Unlike PCA
- No sense of order

## ICA Issues

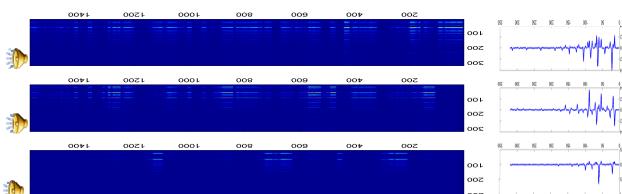
- But the issues here?
- But not much
- Better..



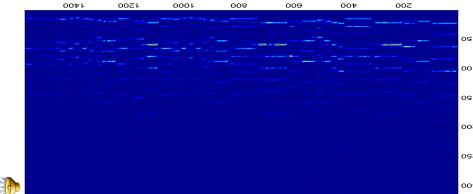
**So how does this work: ICA solution**

- There are 12 notes in the segment, hence we try to estimate 12 notes..

- There are 12 notes in the segment, hence we try to estimate 12 notes..



## PCA Solution



**So how does that work?**