

Machine Learning for Signal Processing

Independent Component Analysis

Class 9. 30 Sep 2014

Instructor: Bhiksha Raj

Correlation vs. Causation

- The consumption of burgers has gone up steadily in the past decade



- In the same period, the penguin population of Antarctica has gone down

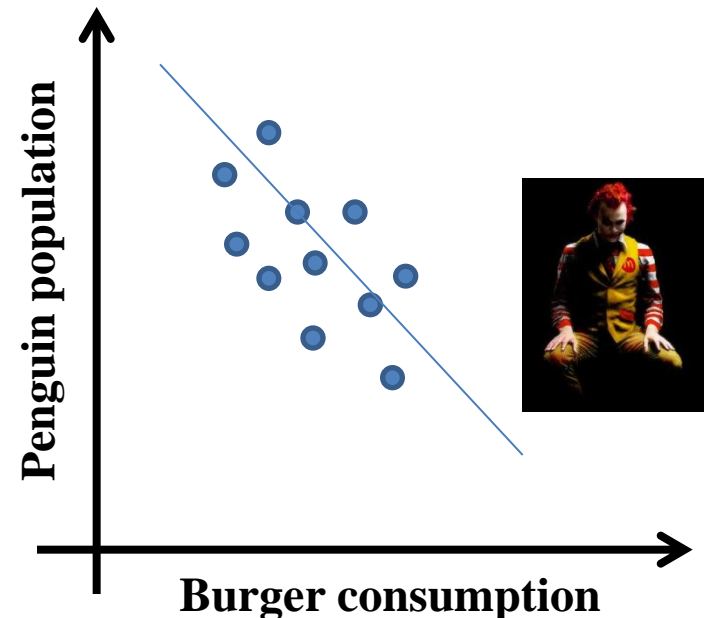
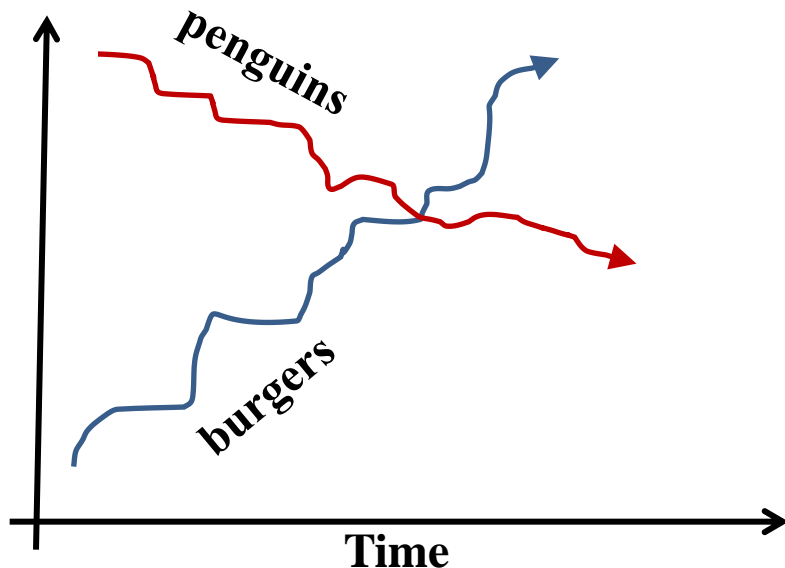


Correlation, not Causation
(unless McDonalds has a
top-secret Antarctica division)



The concept of *correlation*

- Two variables are correlated if knowing the value of one gives you information about the ***expected value*** of the other



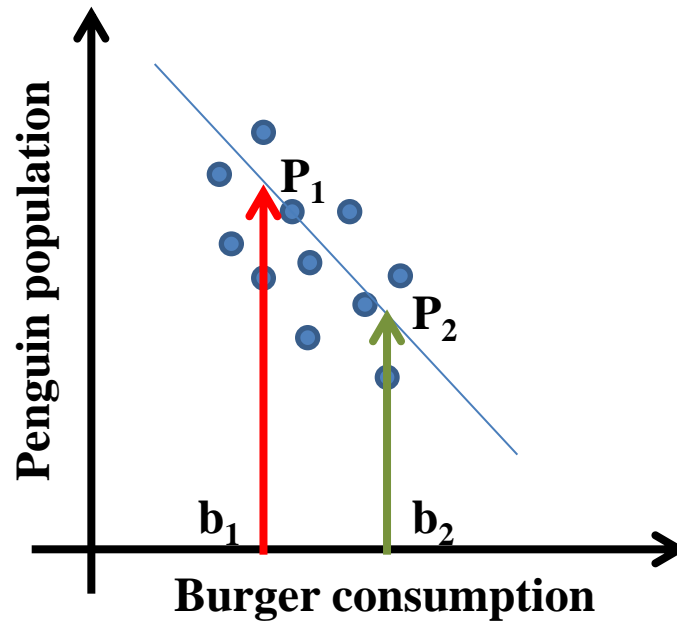
The statistical concept of correlatedness

- Two variables X and Y are correlated if knowing X gives you an *expected* value of Y
- X and Y are uncorrelated if knowing X tells you nothing about the *expected* value of Y
 - Although it could give you other information
 - How?

A brief review of basic probability

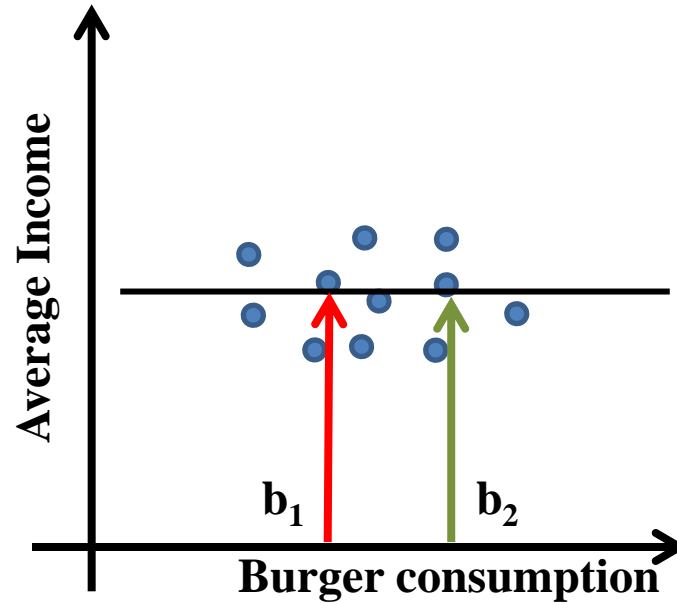
- *Uncorrelated*: Two random variables X and Y are uncorrelated iff:
 - The *average* value of the product of the variables equals the product of their individual averages
- Setup: Each draw produces one instance of X and one instance of Y
 - I.e one instance of (X,Y)
- $E[XY] = E[X]E[Y]$
- The average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

Correlated Variables



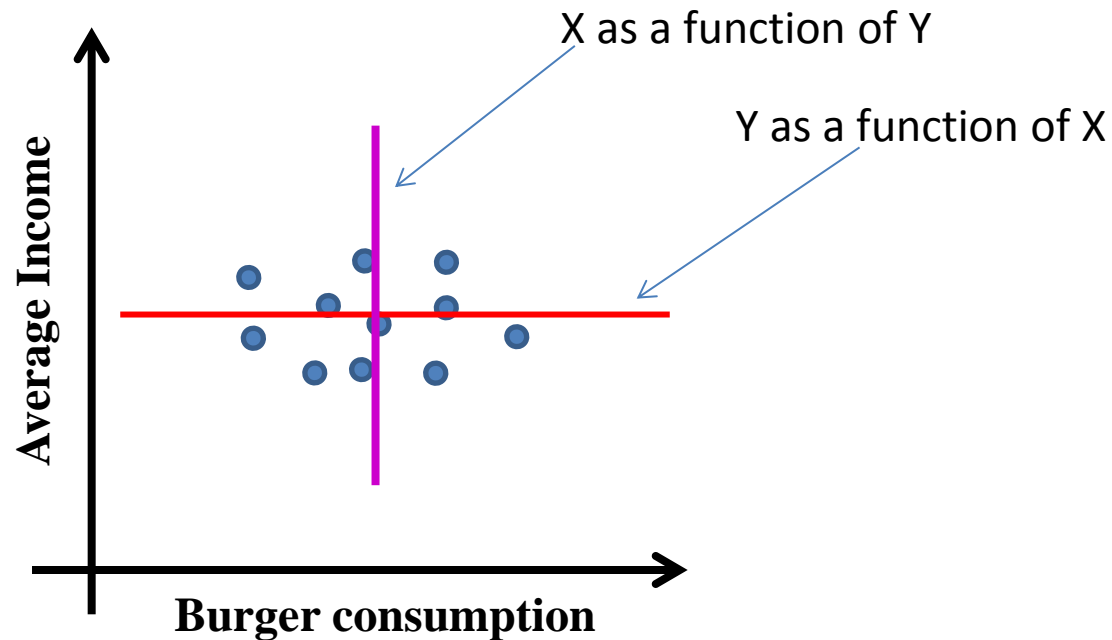
- Expected value of Y given X :
 - Find average of Y values of all samples at (or close) to the given X
 - If this is a function of X , X and Y are correlated

Uncorrelatedness



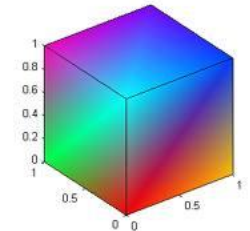
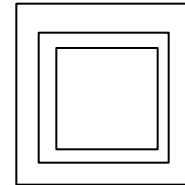
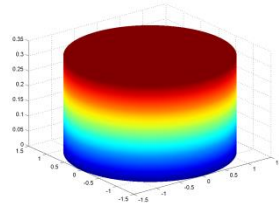
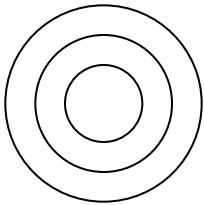
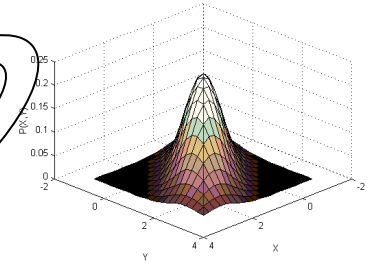
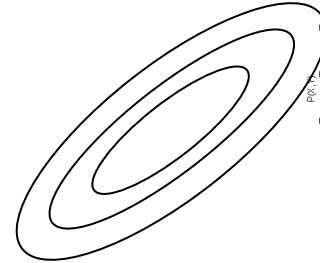
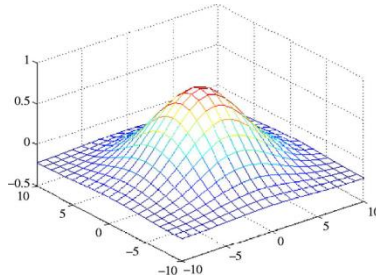
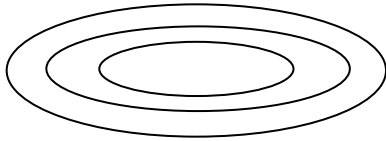
- Knowing X does not tell you what the *average* value of Y is
 - And vice versa

Uncorrelated Variables



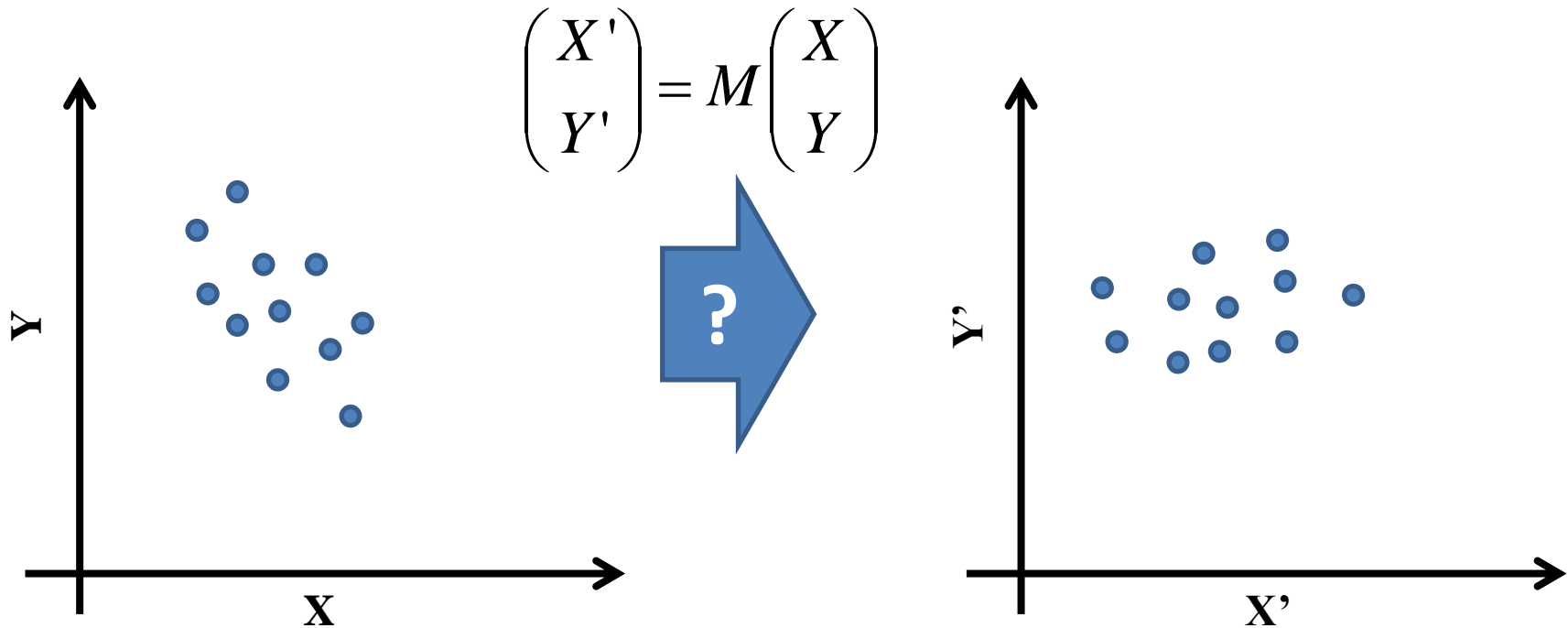
- The average value of Y is the same regardless of the value of X and vice versa

Uncorrelatedness



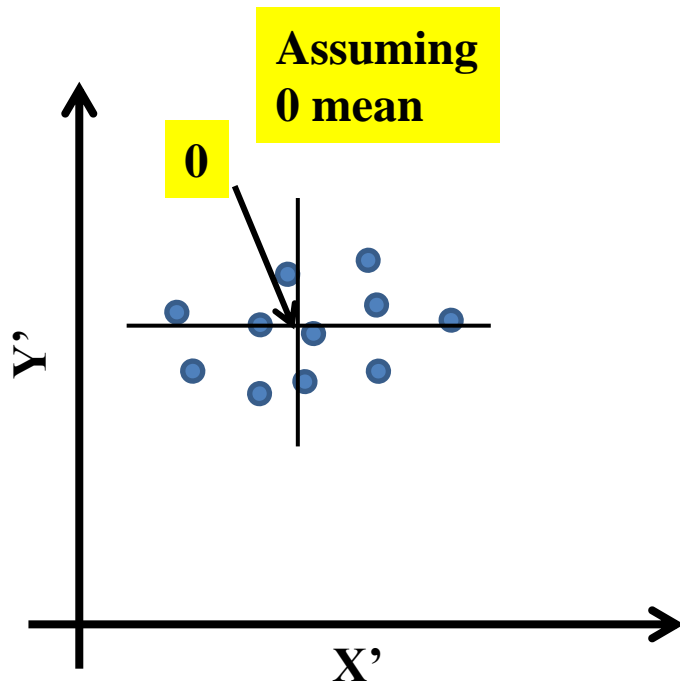
- Which of the above represent uncorrelated RVs?

The notion of *decorrelation*



- So how does one transform the correlated variables (X, Y) to the uncorrelated (X', Y')

What does “decorrelated” mean



- $E[X'] = \text{constant } (0)$
- $E[Y'] = \text{constant } (0)$
- $E[X' | Y'] = 0$
- $E[X'Y'] = E_{Y'}[E[X' | Y']] = 0$

$$E \left[\begin{pmatrix} X' \\ Y' \end{pmatrix} (X' \ Y') \right] = E \begin{pmatrix} X'^2 & X'Y' \\ X'Y' & Y'^2 \end{pmatrix} = \begin{pmatrix} E[X'^2] & 0 \\ 0 & E[Y'^2] \end{pmatrix} = \textit{diagonal matrix}$$

- If \mathbf{Y} is a matrix of vectors, $\mathbf{Y}\mathbf{Y}^T = \text{diagonal}$

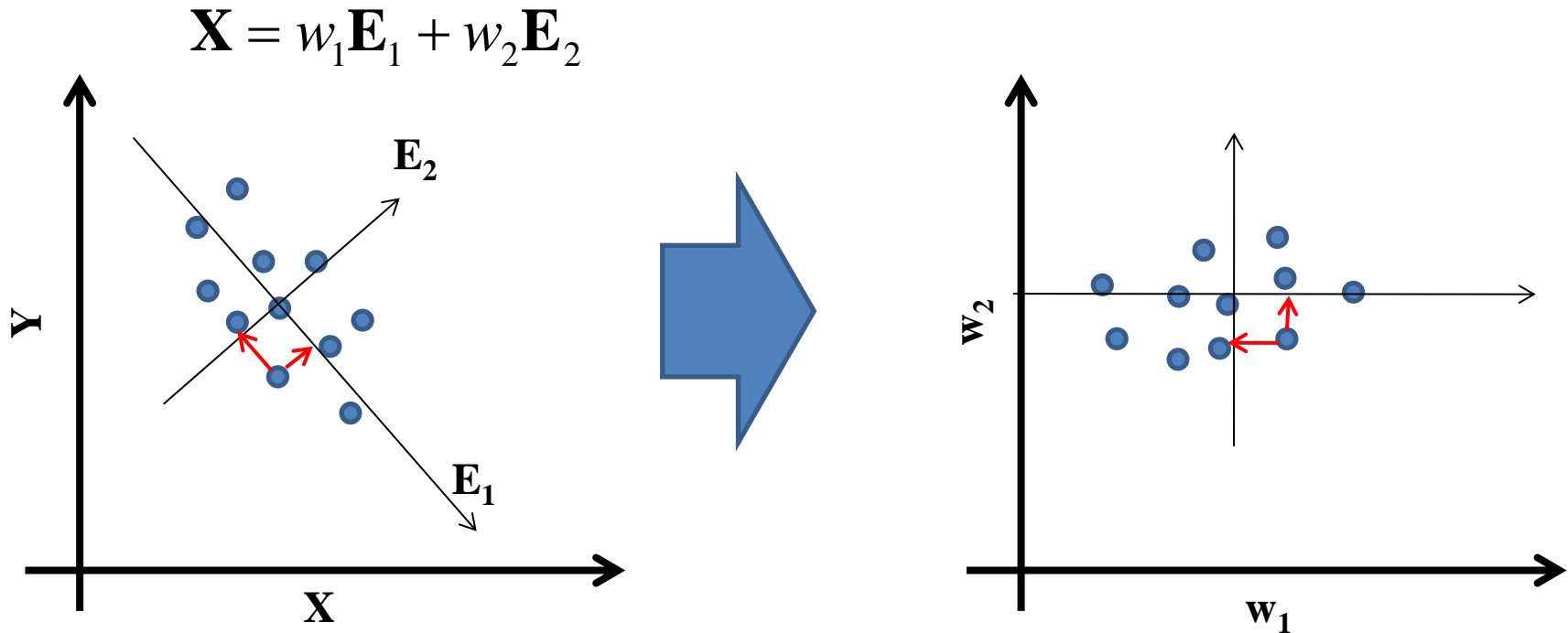
Decorrelation

- Let \mathbf{X} be the matrix of correlated data vectors
 - Each component of \mathbf{X} informs us of the mean trend of other components
- Need a transform \mathbf{M} such that if $\mathbf{Y} = \mathbf{MX}$
- The covariance of \mathbf{Y} is diagonal
 - \mathbf{YY}^T is the covariance if \mathbf{Y} is zero mean
 - $\mathbf{YY}^T = \text{diagonal}$
 - $\Rightarrow \mathbf{MXX}^T\mathbf{M}^T = \mathbf{D}$
 - $\Rightarrow \mathbf{M} \cdot \text{Cov}(\mathbf{X}) \cdot \mathbf{M}^T = \mathbf{D}$

Decorrelation

- Easy solution:
 - Eigen decomposition of $\text{Cov}(\mathbf{X})$: $\text{Cov}(\mathbf{X}) = \mathbf{E}\mathbf{\Lambda}\mathbf{E}^T$
 - $\mathbf{E}\mathbf{E}^T = \mathbf{I}$
- Let $\mathbf{M} = \mathbf{E}^T$
- $\mathbf{M}\text{Cov}(\mathbf{X})\mathbf{M}^T = \mathbf{E}^T\mathbf{E}\mathbf{\Lambda}\mathbf{E}^T\mathbf{E} = \mathbf{\Lambda} = \text{diagonal}$
- PCA: $\mathbf{Y} = \mathbf{M}\mathbf{X}$
- *Diagonalizes* the covariance matrix
 - “Decorrelates” the data

PCA



- PCA: $\mathbf{Y} = \mathbf{M}\mathbf{X}$
- *Diagonalizes* the covariance matrix
 - “Decorrelates” the data

The statistical concept of *Independence*

- Two variables X and Y are *dependent* if knowing X gives you *any information about* Y
- X and Y are *independent* if knowing X tells you nothing at all of Y

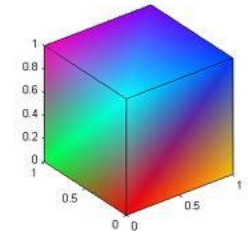
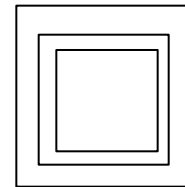
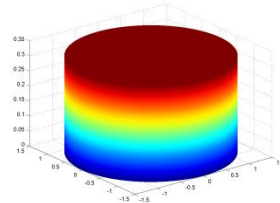
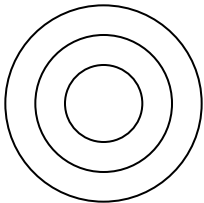
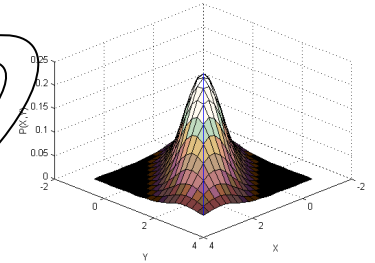
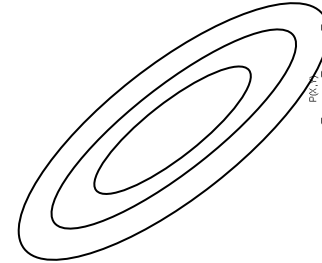
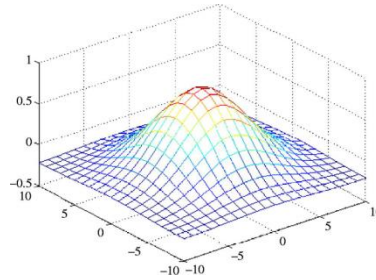
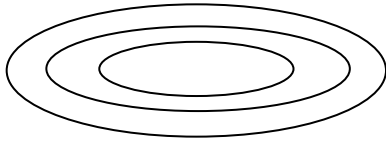
A brief review of basic probability

- ***Independence***: Two random variables X and Y are independent iff:
 - Their joint probability equals the product of their individual probabilities
- $P(X,Y) = P(X)P(Y)$
- Independence implies uncorrelatedness
 - The average value of X is the same regardless of the value of Y
 - $E[X|Y] = E[X]$
 - But not the other way

A brief review of basic probability

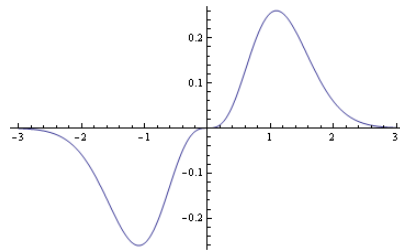
- *Independence*: Two random variables X and Y are independent iff:
- The average value of *any function* of X is the same regardless of the value of Y
 - Or any function of Y
- $E[f(X)g(Y)] = E[f(X)] E[g(Y)]$ for all $f()$, $g()$

Independence

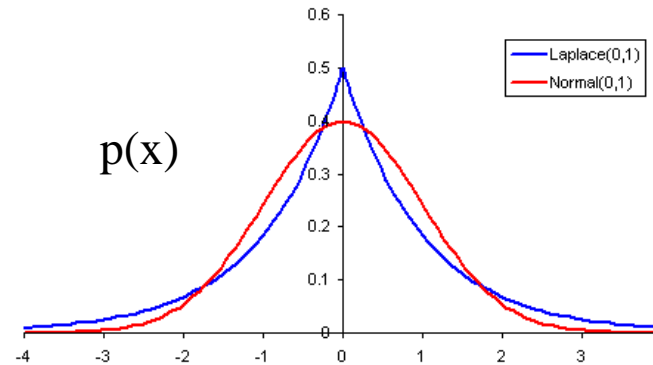
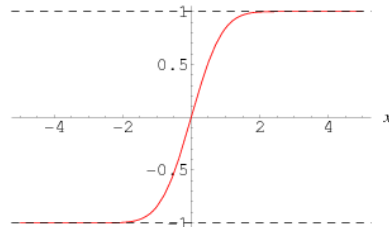


- Which of the above represent independent RVs?
- Which represent uncorrelated RVs?

A brief review of basic probability



$$y = f(x)$$



- The expected value of an odd function of an RV is 0 if
 - The RV is 0 mean
 - The PDF of the RV is symmetric around 0
- **$E[f(X)] = 0$ if $f(X)$ is odd symmetric**

A brief review of basic info. theory



T(all), M(ed), S(hort)...

$$H(X) = \sum_X P(X) [-\log P(X)]$$

- Entropy: The *minimum average* number of bits to transmit to convey a symbol



T, M, S...

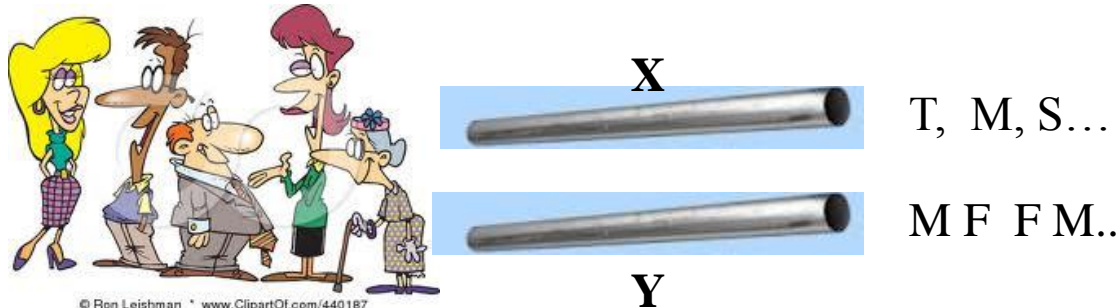


M F F M..

$$H(X,Y) = \sum_{X,Y} P(X,Y) [-\log P(X,Y)]$$

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

A brief review of basic info. theory



$$H(X | Y) = \sum_Y P(Y) \sum_X P(X | Y) [-\log P(X | Y)] = \sum_{X,Y} P(X, Y) [-\log P(X | Y)]$$

- Conditional Entropy: The *minimum average* number of bits to transmit to convey a symbol X , after symbol Y has already been conveyed
 - Averaged over all values of X and Y

A brief review of basic info. theory

$$H(X | Y) = \sum_Y P(Y) \sum_X P(X | Y) [-\log P(X | Y)] = \sum_Y P(Y) \sum_X P(X) [-\log P(X)] = H(X)$$

- Conditional entropy of $X = H(X)$ if X is independent of Y

$$H(X, Y) = \sum_{X, Y} P(X, Y) [-\log P(X, Y)] = \sum_{X, Y} P(X, Y) [-\log P(X)P(Y)]$$

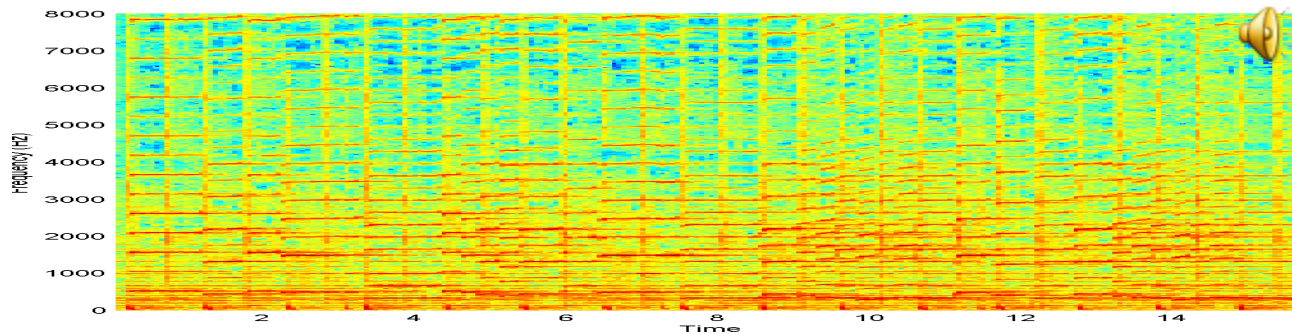
$$= -\sum_{X, Y} P(X, Y) \log P(X) - \sum_{X, Y} P(X, Y) \log P(Y) = H(X) + H(Y)$$

- Joint entropy of X and Y is the sum of the entropies of X and Y if they are independent

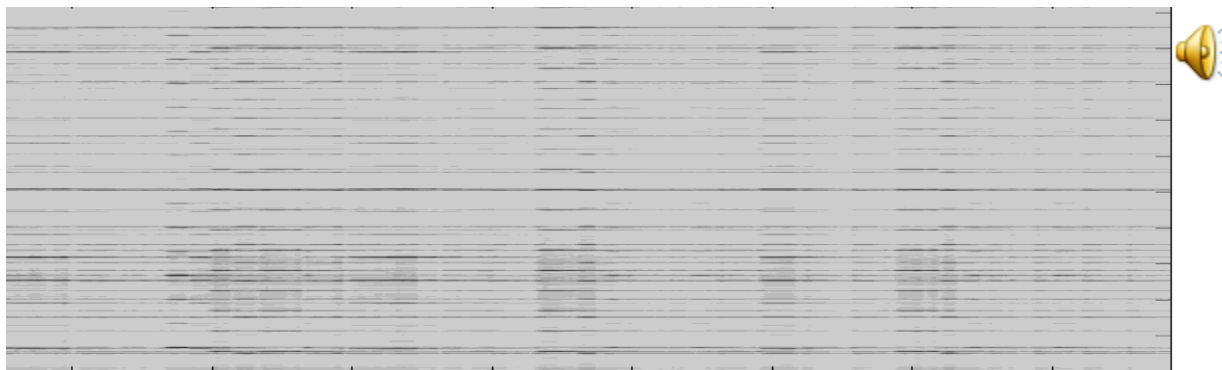
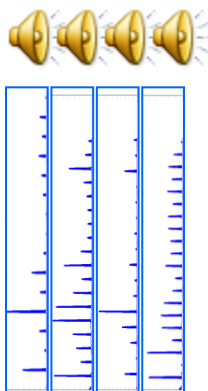
Onward..

Projection: multiple notes

M =



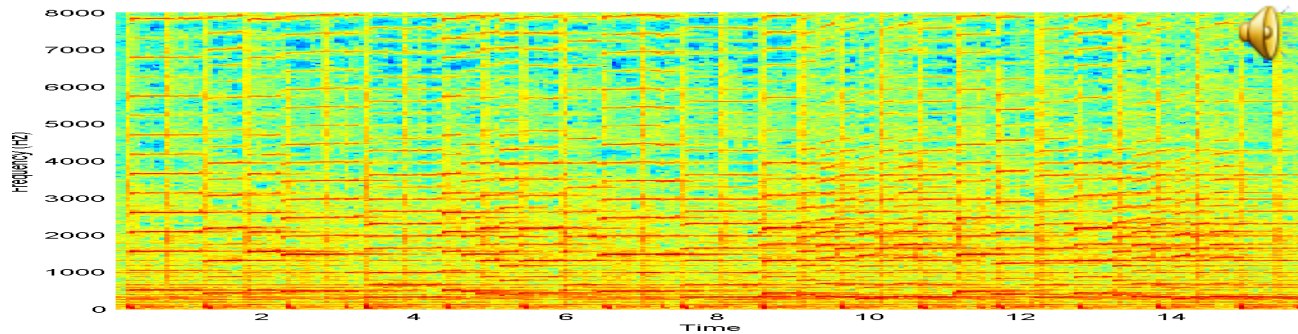
W =



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

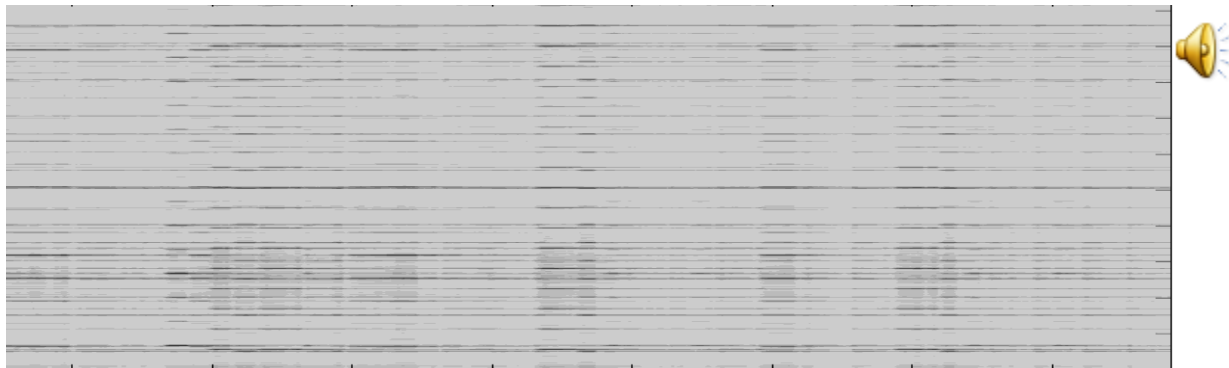
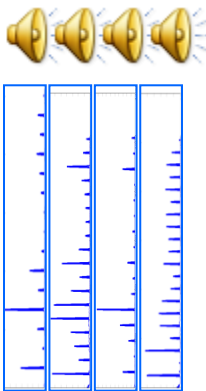
We're actually computing a score

M =



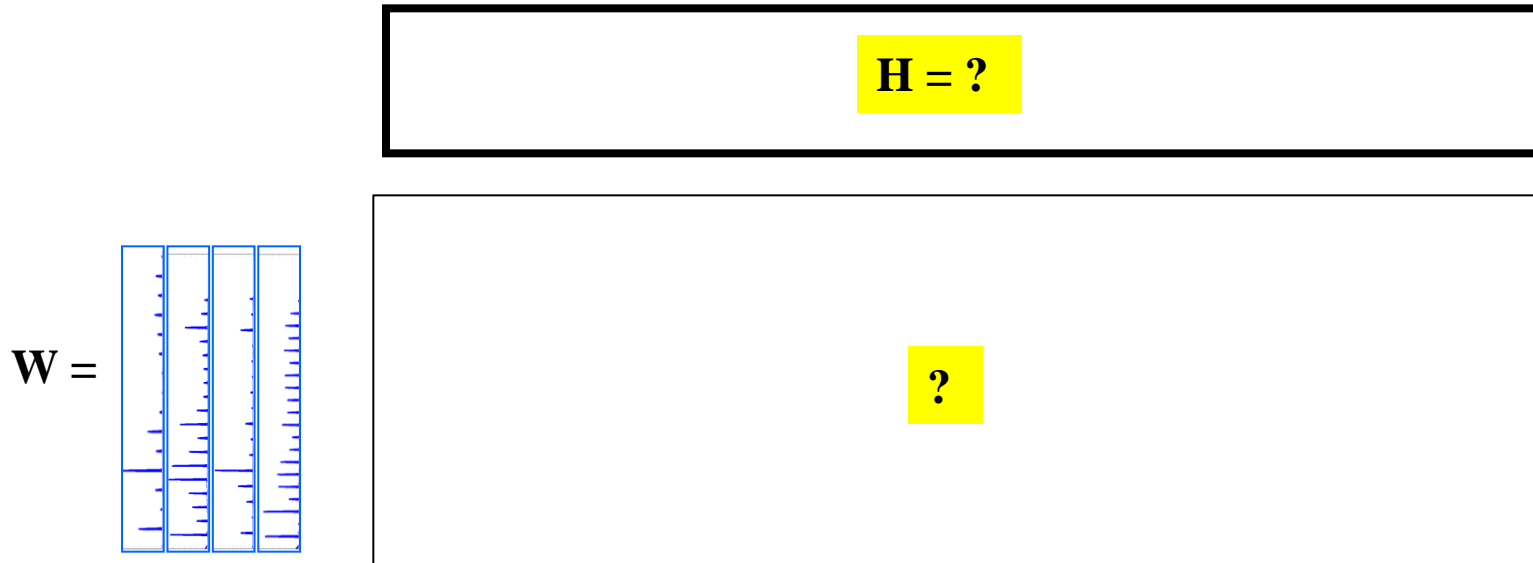
H = ?

W =



- $M \sim WH$
- $H = \text{pinv}(W)M$

So what are we doing here?

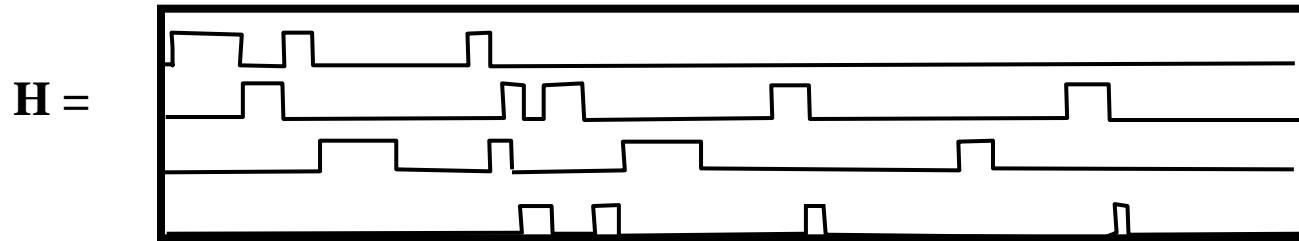
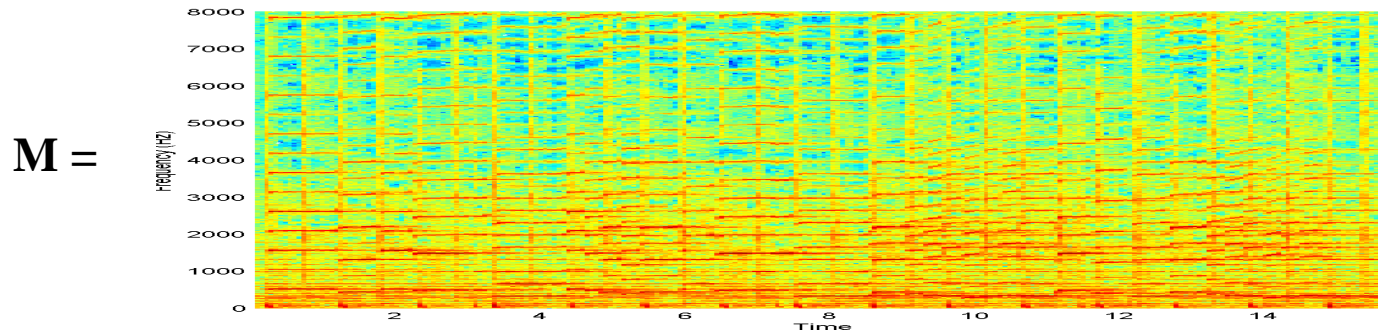


- $\mathbf{M} \sim \mathbf{WH}$ is an approximation
- Given \mathbf{W} , estimate \mathbf{H} to minimize error

$$\mathbf{H} = \arg \min_{\bar{\mathbf{H}}} \|\mathbf{M} - \mathbf{W}\bar{\mathbf{H}}\|_F^2 = \arg \min_{\bar{\mathbf{H}}} \sum_i \sum_j (\mathbf{M}_{ij} - (\mathbf{W}\bar{\mathbf{H}})_{ij})^2$$

- Must ideally find *transcription* of given notes

How about the other way?



W = ?

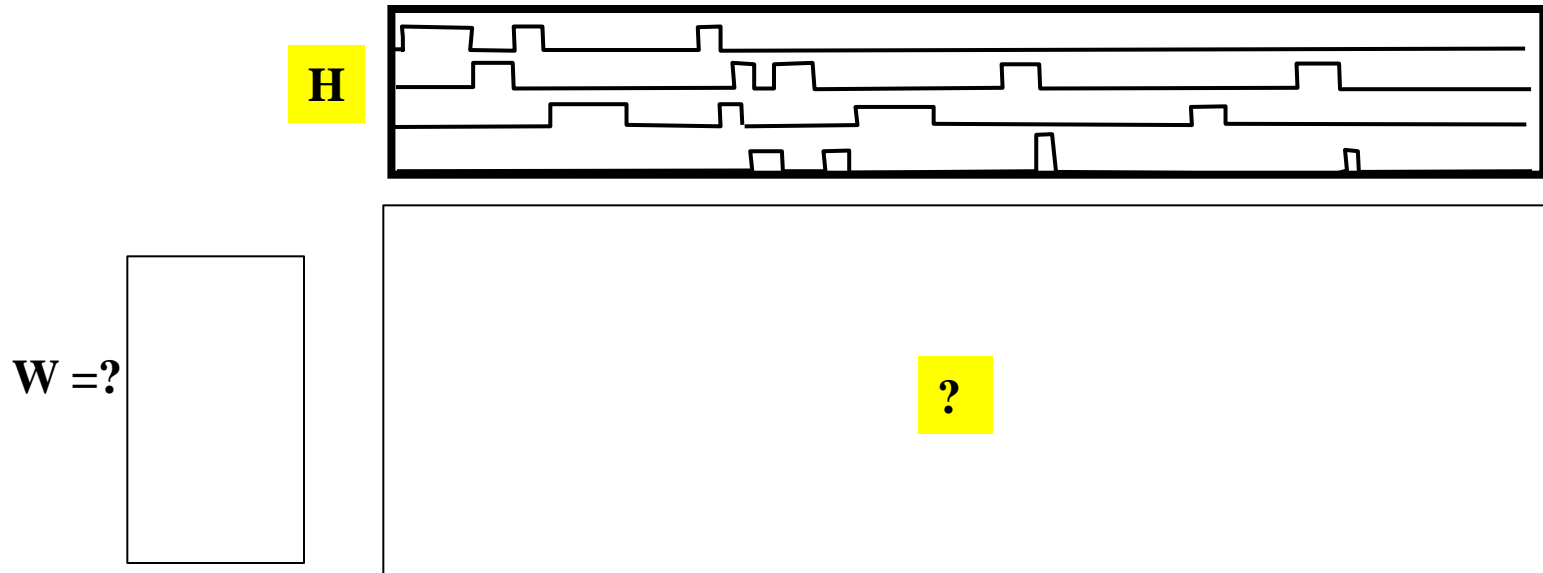
U = ?

■ $M \sim WH$

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

$$U = WH$$

Going the other way..



- $\mathbf{M} \sim \mathbf{W}\mathbf{H}$ is an approximation
- Given \mathbf{H} , estimate \mathbf{W} to minimize error

$$\mathbf{W} = \arg \min_{\bar{\mathbf{W}}} \|\mathbf{M} - \bar{\mathbf{W}}\mathbf{H}\|_F^2 = \arg \min_{\bar{\mathbf{H}}} \sum_i \sum_j (\mathbf{M}_{ij} - (\bar{\mathbf{W}}\mathbf{H})_{ij})^2$$

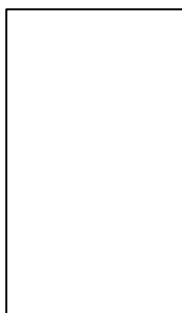
- Must ideally find the *notes* corresponding to the transcription

When both parameters are unknown

H = ?



W = ?



approx(M) = ?

- Must estimate both **H** and **W** to best approximate **M**
- Ideally, must learn *both* the *notes* and *their* transcription!

A least squares solution

$$\mathbf{W}, \mathbf{H} = \arg \min_{\overline{\mathbf{W}}, \overline{\mathbf{H}}} \|\mathbf{M} - \overline{\mathbf{W}}\overline{\mathbf{H}}\|_F^2$$

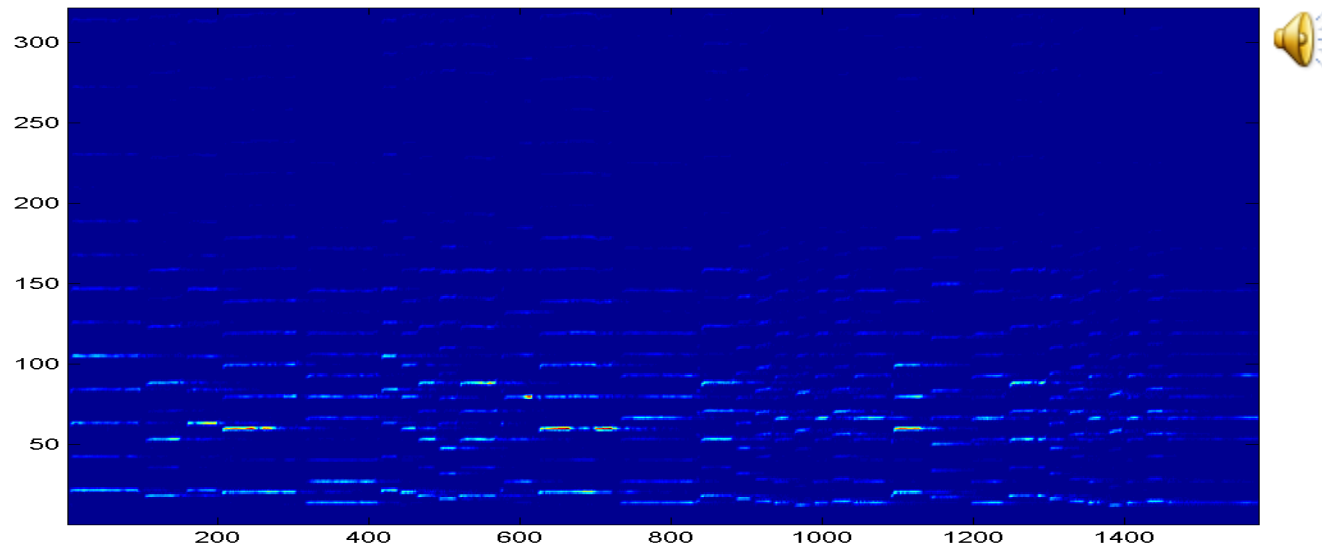
- Unconstrained
 - For any \mathbf{W}, \mathbf{H} that minimizes the error, $\mathbf{W}' = \mathbf{W}\mathbf{A}$, $\mathbf{H}' = \mathbf{A}^{-1}\mathbf{H}$ also minimizes the error for any invertible \mathbf{A}
 - Too many solutions
- Constraint: \mathbf{W} is orthogonal
 - $\mathbf{W}^T\mathbf{W} = \mathbf{I}$
 - PCA!!

PCA: Constrained solution

$$\mathbf{W}, \mathbf{H} = \arg \min_{\overline{\mathbf{W}}, \overline{\mathbf{H}}} \|\mathbf{M} - \overline{\mathbf{W}}\overline{\mathbf{H}}\|_F^2$$

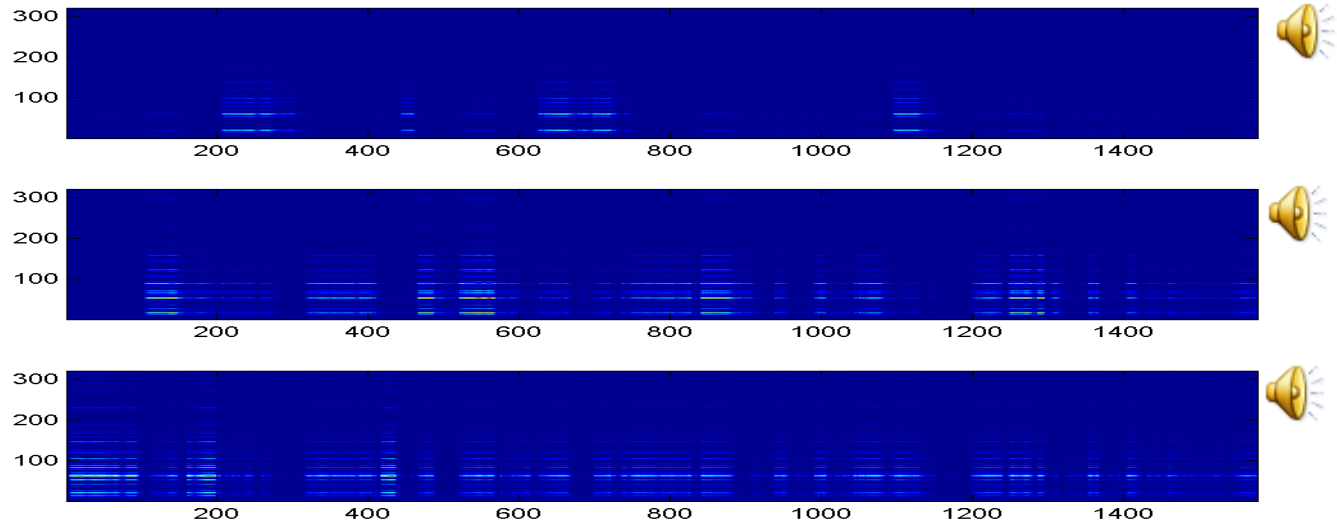
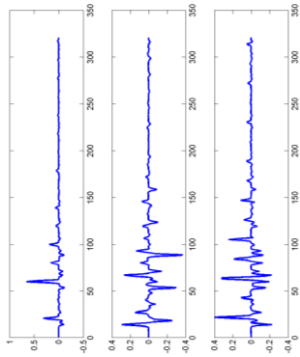
- Constraint: \mathbf{W} is orthogonal
 - $\mathbf{W}^T \mathbf{W} = \mathbf{I}$
- This results in PCA!!
 - \mathbf{W} are the Eigenvectors of $\mathbf{M}\mathbf{M}^T$

So how does that work?



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

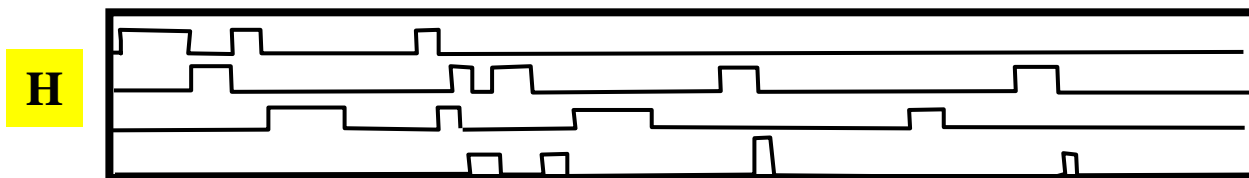
So how does that work?



- There are 12 notes in the segment, hence we try to estimate 12 notes..
- Results are not good

A *constrained* least squares solution

$$\mathbf{W}, \mathbf{H} = \arg \min_{\overline{\mathbf{W}}, \overline{\mathbf{H}}} \|\mathbf{M} - \overline{\mathbf{W}}\overline{\mathbf{H}}\|_F^2$$



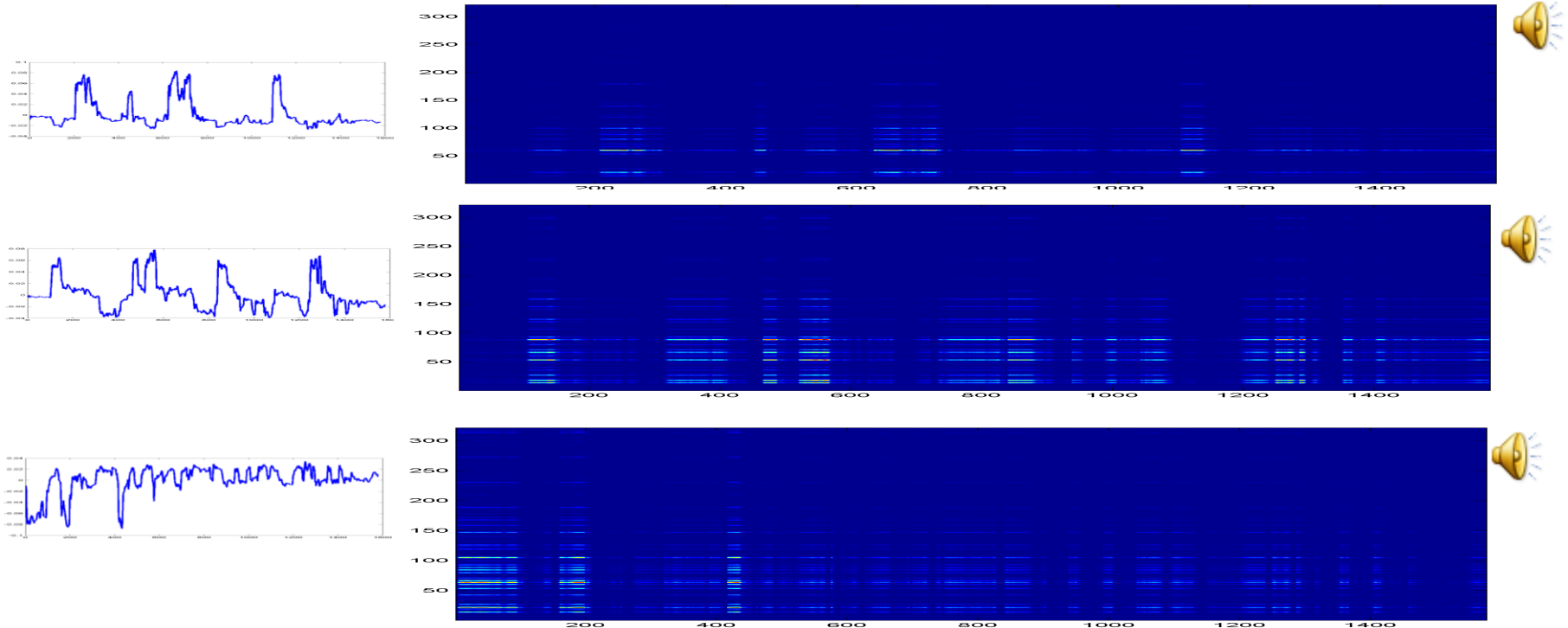
- For our problem, let's consider the “truth”..
- When one note occurs, the other does not
 - $\mathbf{h}_i^T \mathbf{h}_j = 0$ for all $i \neq j$
- The rows of \mathbf{H} are *uncorrelated*

PCA: The *Other* Way?

$$\mathbf{W}, \mathbf{H} = \arg \min_{\overline{\mathbf{W}}, \overline{\mathbf{H}}} \|\mathbf{M} - \overline{\mathbf{W}}\overline{\mathbf{H}}\|_F^2$$

- Constraint: \mathbf{H} is orthogonal
 - $\mathbf{H}\mathbf{H}^T = \mathbf{I}$
- This results in PCA or the row vectors of \mathbf{M} !!
 - \mathbf{H} are the Eigenvectors of $\mathbf{M}^T\mathbf{M}$

So how does that work?



- The scores of the first three “notes” and their contributions
- Not that great again

PCA

H = ?

W = ?

approx(M) = ?

- If the notes matrix **W** is made orthogonal, the rows of **H** end up being orthogonal to one another
 - **H** is the orthogonalized version of **M**
- If the scores matrix **H** is made orthogonal instead, the rows of **W** end up being orthogonal
- The two decompositions are identical to within a scaling of the vectors

Eigendecomposition and SVD

$$\mathbf{M} = \mathbf{U}\mathbf{S}\mathbf{V}^T$$

$$\mathbf{M} = \mathbf{W}\mathbf{H}$$

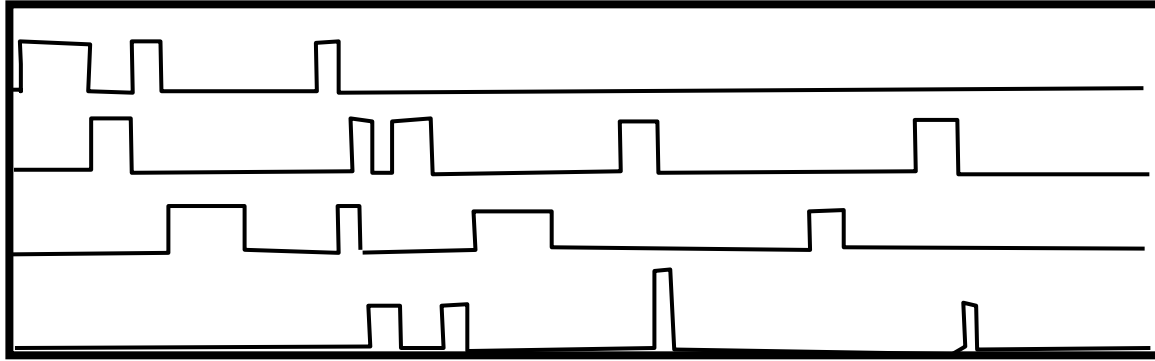
- Matrix \mathbf{M} can be decomposed as $\mathbf{M} = \mathbf{U}\mathbf{S}\mathbf{V}^T$
- When we assume the scores are orthogonal, we get
$$\mathbf{H} = \mathbf{V}^T, \mathbf{W} = \mathbf{U}\mathbf{S}$$
- When we assume the notes are orthogonal, we get
$$\mathbf{W} = \mathbf{U}, \mathbf{H} = \mathbf{S}\mathbf{V}^T$$
- **In either case the results are the same**
 - The notes are orthogonal and so are the scores
 - Not good in our problem

Orthogonality

$$\mathbf{M} = \mathbf{W}\mathbf{H}$$

- In any *least-squared error* decomposition $\mathbf{M} = \mathbf{W}\mathbf{H}$, if the columns of \mathbf{W} are orthogonal, the rows of \mathbf{H} will also be orthogonal
- Sometimes mere orthogonality is not enough

What *else* can we look for?



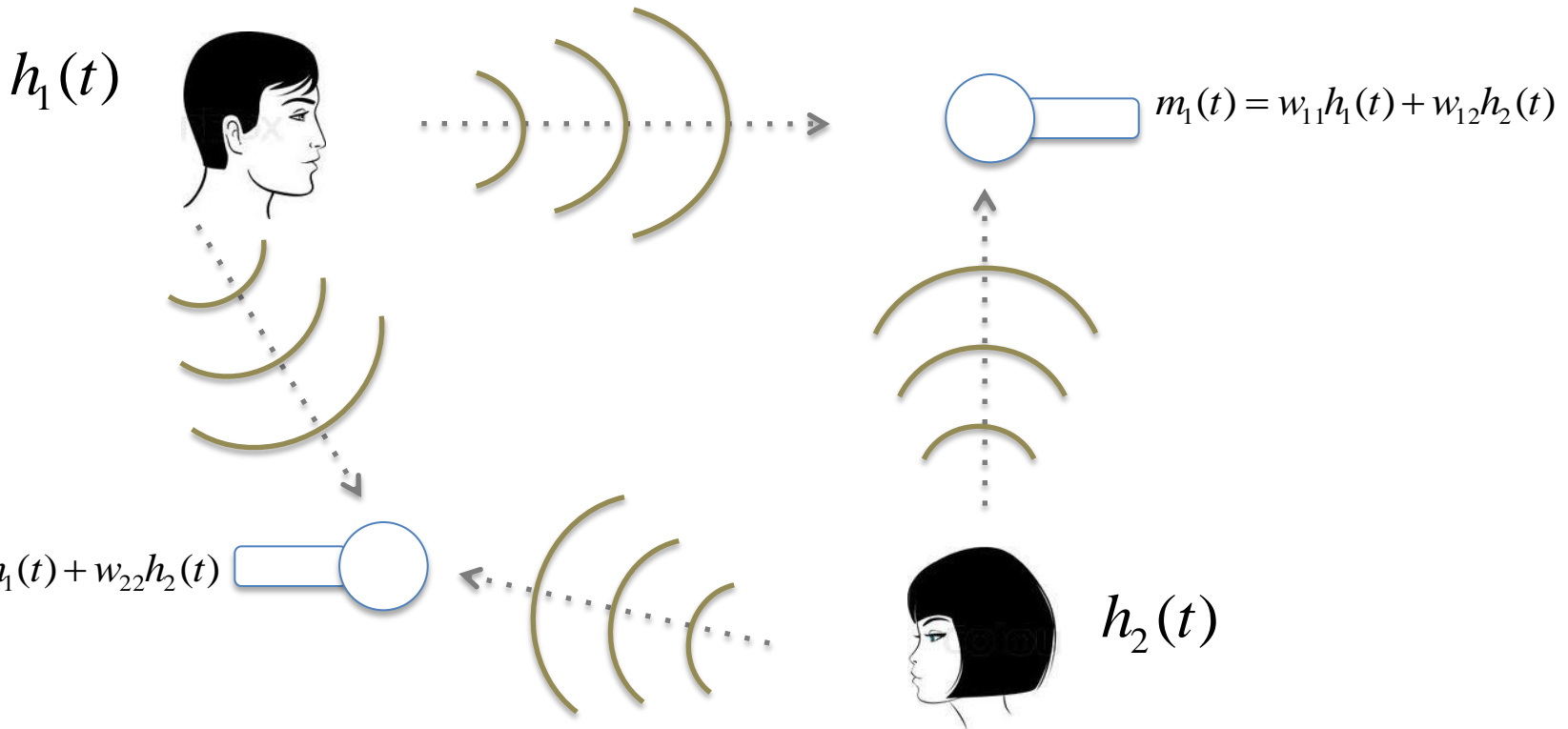
- Assume: The “transcription” of one note does not depend on what else is playing
 - Or, in a multi-instrument piece, instruments are playing independently of one another
- Not strictly true, but still..

Formulating it with Independence

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

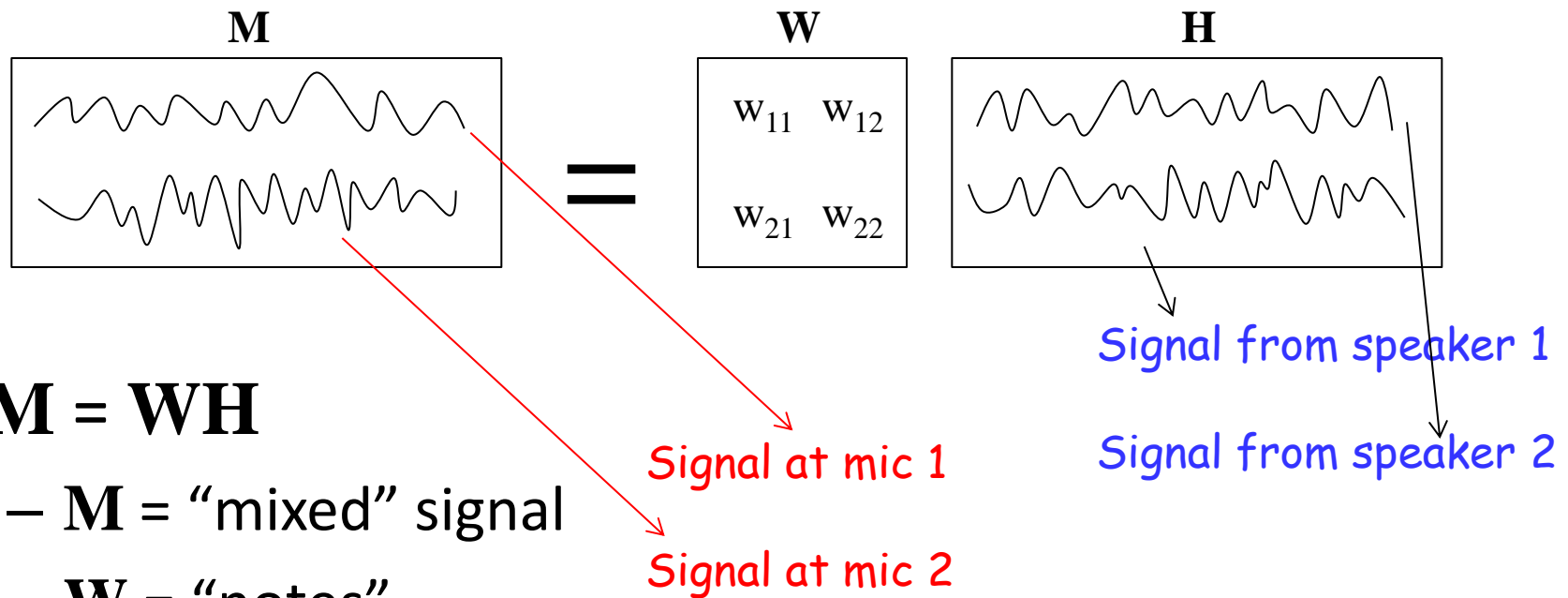
- Impose statistical independence constraints on decomposition

Changing problems for a bit



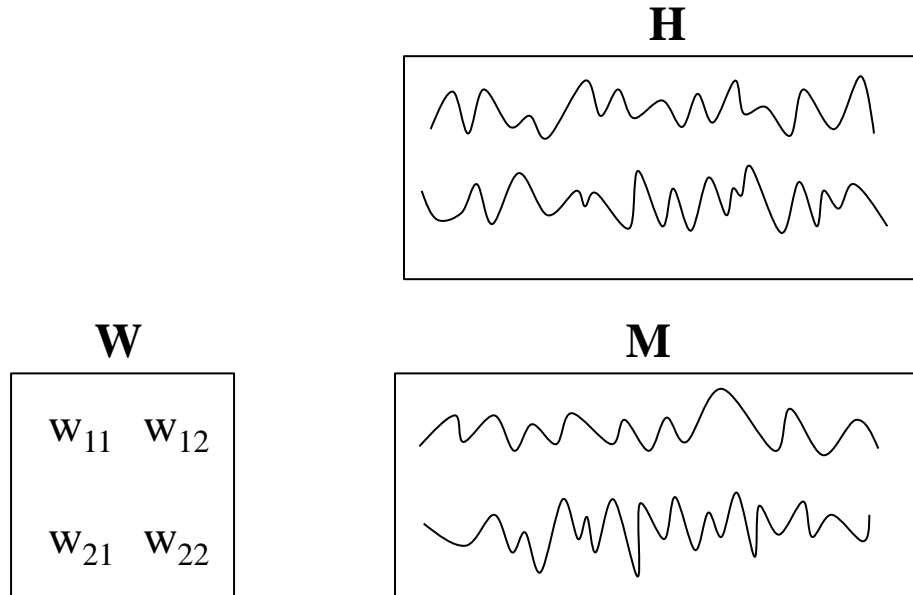
- Two people speak simultaneously
- Recorded by two microphones
- Each recorded signal is a mixture of both signals

A Separation Problem



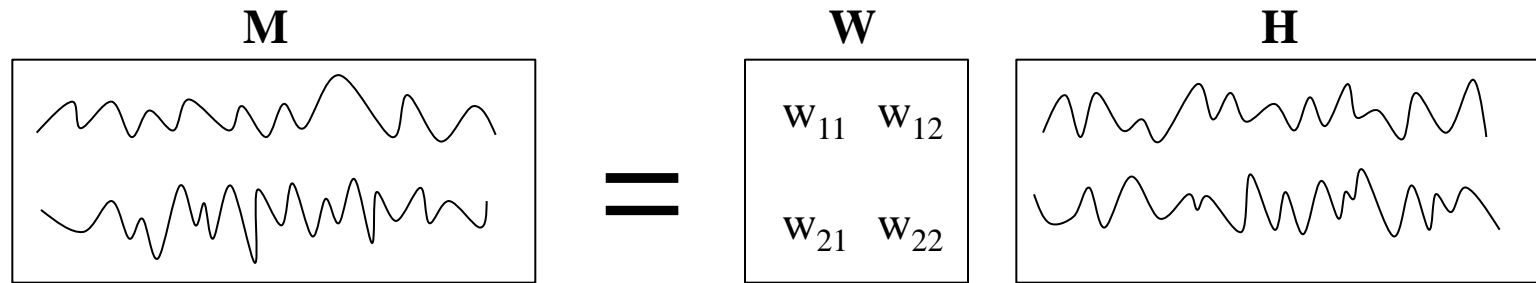
- $\mathbf{M} = \mathbf{W}\mathbf{H}$
 - \mathbf{M} = “mixed” signal
 - \mathbf{W} = “notes”
 - \mathbf{H} = “transcription”
- Separation challenge: Given only \mathbf{M} estimate \mathbf{H}
- Identical to the problem of “finding notes”

A Separation Problem



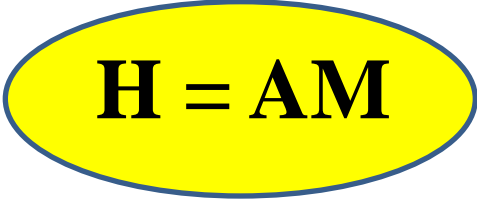
- Separation challenge: Given only **M** estimate **H**
- **Identical to the problem of “finding notes”**

Imposing Statistical Constraints

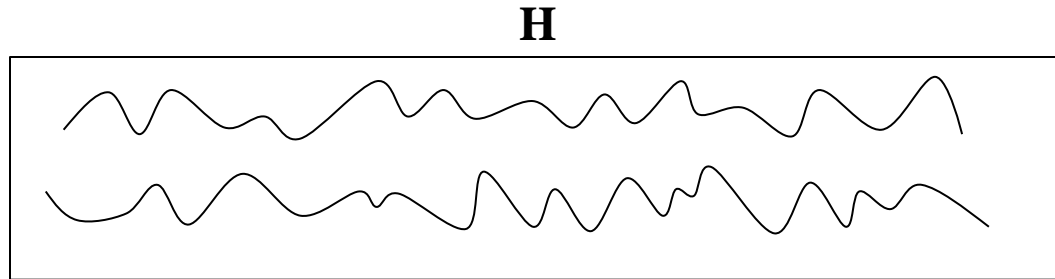


- $\mathbf{M} = \mathbf{W}\mathbf{H}$
- Given only \mathbf{M} estimate \mathbf{H}
- $\mathbf{H} = \mathbf{W}^{-1}\mathbf{M} = \mathbf{A}\mathbf{M}$
- Only known constraint: The rows of \mathbf{H} are independent
- Estimate \mathbf{A} such that the components of $\mathbf{A}\mathbf{M}$ are statistically independent
 - \mathbf{A} is the *unmixing* matrix

Statistical Independence

- $\mathbf{M} = \mathbf{W}\mathbf{H}$ $\mathbf{H} = \mathbf{A}\mathbf{M}$  Remember this form
- *Emulating independence*
 - Compute \mathbf{W} (or \mathbf{A}) and \mathbf{H} such that \mathbf{H} has statistical characteristics that are observed in statistically independent variables
- *Enforcing independence*
 - Compute \mathbf{W} and \mathbf{H} such that the components of \mathbf{M} are independent

Emulating Independence



- The rows of **H** are uncorrelated
 - $E[\mathbf{h}_i \mathbf{h}_j] = E[\mathbf{h}_i]E[\mathbf{h}_j]$
 - \mathbf{h}_i and \mathbf{h}_j are the i^{th} and j^{th} components of any vector in **H**
- The fourth order moments are independent
 - $E[\mathbf{h}_i \mathbf{h}_j \mathbf{h}_k \mathbf{h}_l] = E[\mathbf{h}_i]E[\mathbf{h}_j]E[\mathbf{h}_k]E[\mathbf{h}_l]$
 - $E[\mathbf{h}_i^2 \mathbf{h}_j \mathbf{h}_k] = E[\mathbf{h}_i^2]E[\mathbf{h}_j]E[\mathbf{h}_k]$
 - $E[\mathbf{h}_i^2 \mathbf{h}_j^2] = E[\mathbf{h}_i^2]E[\mathbf{h}_j^2]$
 - Etc.

Zero Mean

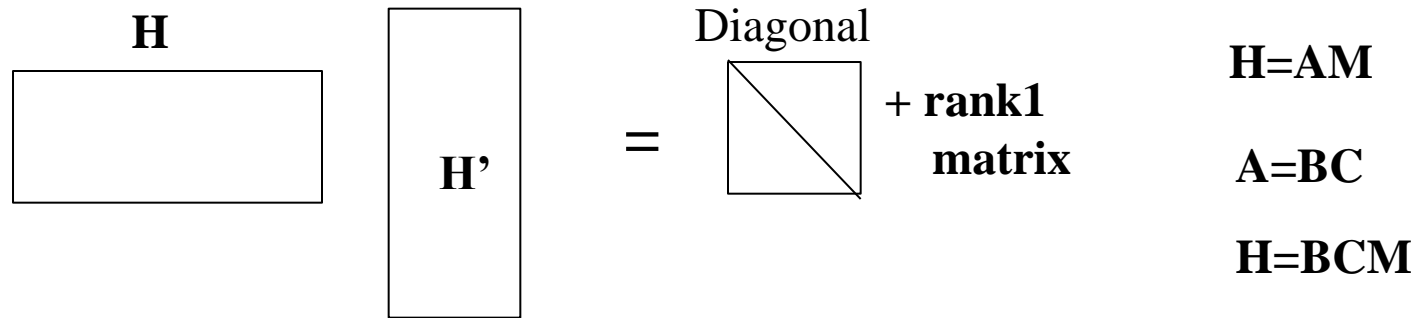
- Usual to assume *zero mean* processes
 - Otherwise, some of the math doesn't work well
- $\mathbf{M} = \mathbf{W}\mathbf{H}$ $\mathbf{H} = \mathbf{A}\mathbf{M}$
- If $\text{mean}(\mathbf{M}) = \mathbf{0} \Rightarrow \text{mean}(\mathbf{H}) = \mathbf{0}$
 - $\mathbf{E}[\mathbf{H}] = \mathbf{A} \cdot \mathbf{E}[\mathbf{M}] = \mathbf{A}\mathbf{0} = \mathbf{0}$
 - First step of ICA: Set the mean of \mathbf{M} to $\mathbf{0}$

$$\mu_{\mathbf{m}} = \frac{1}{\text{cols}(\mathbf{M})} \sum_i \mathbf{m}_i$$

$$\mathbf{m}_i = \mathbf{m}_i - \mu_{\mathbf{m}} \quad \forall i$$

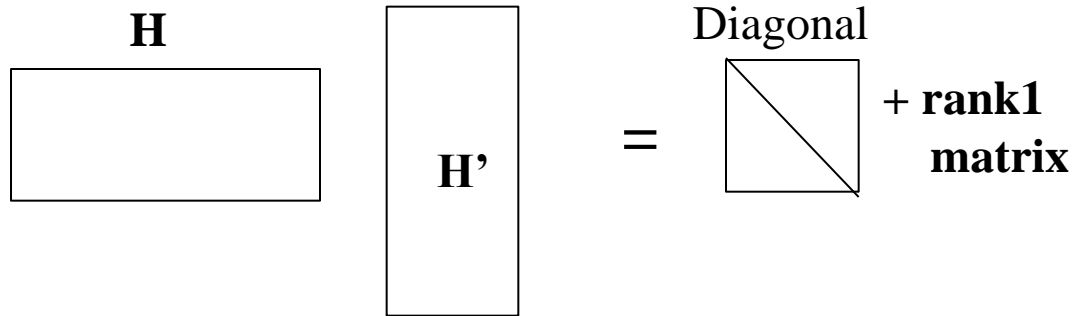
- \mathbf{m}_i are the columns of \mathbf{M}

Emulating Independence..



- Independence \rightarrow Uncorrelatedness
- Estimate a **C** such that **CM** is uncorrelated
- **X = CM**
 - $E[\mathbf{x}_i \mathbf{x}_j] = \delta_{ij}$ [since **M** is now “centered”]
 - **XX^T = I**
 - In reality, we only want this to be a diagonal matrix, but we’ll make it identity

Decorrelating



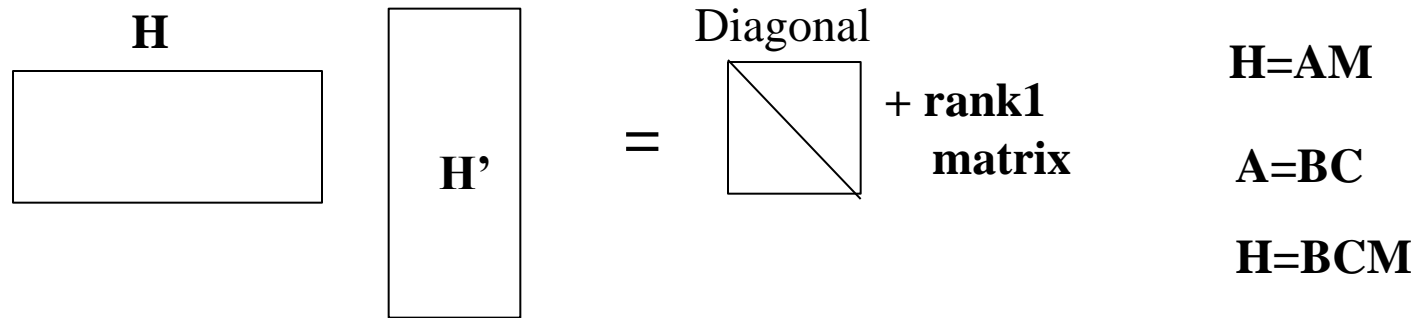
$$\mathbf{H} = \mathbf{A}\mathbf{M}$$

$$\mathbf{A} = \mathbf{B}\mathbf{C}$$

$$\mathbf{H} = \mathbf{B}\mathbf{C}\mathbf{M}$$

- $\mathbf{X} = \mathbf{C}\mathbf{M}$
- $\mathbf{X}\mathbf{X}^T = \mathbf{I}$
- Eigen decomposition $\mathbf{M}\mathbf{M}^T = \mathbf{E}\mathbf{S}\mathbf{E}^T$
- Let $\mathbf{C} = \mathbf{S}^{-1/2}\mathbf{E}^T$
 - $\mathbf{X} = \mathbf{S}^{-1/2}\mathbf{E}^T\mathbf{M}$
 - $\mathbf{X}\mathbf{X}^T = \mathbf{C}\mathbf{M}\mathbf{M}^T\mathbf{C}^T = \mathbf{S}^{-1/2}\mathbf{E}^T\mathbf{E}\mathbf{S}\mathbf{E}^T\mathbf{E}\mathbf{S}^{-1/2} = \mathbf{I}$

Decorrelating



- Eigen decomposition $\mathbf{M}\mathbf{M}^T = \mathbf{E}\mathbf{S}\mathbf{E}^T$
- Let $\mathbf{C} = \mathbf{S}^{-1/2}\mathbf{E}^T$
- $\mathbf{X} = \mathbf{C}\mathbf{M}$

- $\mathbf{X}\mathbf{X}^T = \mathbf{I}$

- \mathbf{X} is called the *whitened* version of \mathbf{M}
 - The process of decorrelating \mathbf{M} is called *whitening*
 - \mathbf{C} is the *whitening matrix*

Uncorrelated != Independent

- Whitening merely ensures that the resulting signals are uncorrelated, i.e.

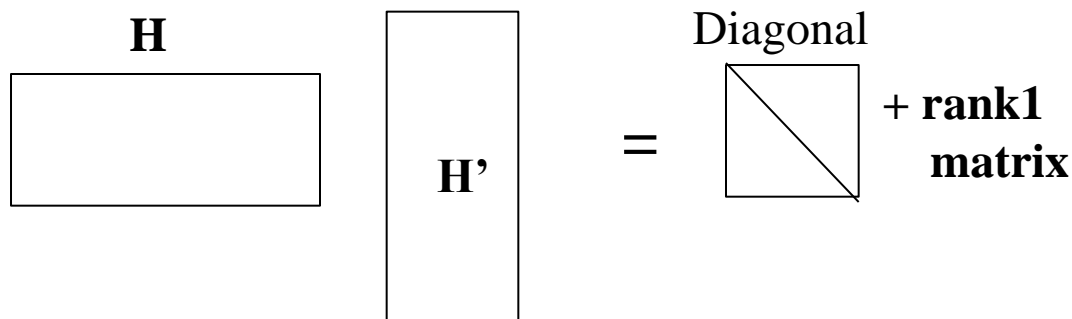
$$E[\mathbf{x}_i \mathbf{x}_j] = 0 \text{ if } i \neq j$$

- This does not ensure higher order moments are also decoupled, e.g. it does not ensure that

$$E[\mathbf{x}_i^2 \mathbf{x}_j^2] = E[\mathbf{x}_i^2] E[\mathbf{x}_j^2]$$

- This is *one* of the signatures of independent RVs
- Lets explicitly decouple the fourth order moments

Decorrelating



$$\mathbf{H} = \mathbf{A}\mathbf{M}$$

$$\mathbf{A} = \mathbf{B}\mathbf{C}$$

$$\mathbf{H} = \mathbf{B}\mathbf{C}\mathbf{M}$$

$$\mathbf{H} = \mathbf{B}\mathbf{X}$$

- $\mathbf{X} = \mathbf{C}\mathbf{M}$
- $\mathbf{X}\mathbf{X}^T = \mathbf{I}$
- Will multiplying \mathbf{X} by \mathbf{B} *re-correlate* the components?
- Not if \mathbf{B} is *unitary*
 - $\mathbf{B}\mathbf{B}^T = \mathbf{B}^T\mathbf{B} = \mathbf{I}$
- $\mathbf{H}\mathbf{H}^T = \mathbf{B}\mathbf{X}\mathbf{X}^T\mathbf{B}^T = \mathbf{B}\mathbf{B}^T = \mathbf{I}$
- So we want to find a *unitary* matrix
 - Since the rows of \mathbf{H} are uncorrelated
 - Because they are independent

ICA: Freeing Fourth Moments

- We have $E[\mathbf{x}_i \mathbf{x}_j] = 0$ if $i \neq j$
 - Already been decorrelated
- $\mathbf{A}=\mathbf{BC}$, $\mathbf{H} = \mathbf{BCM}$, $\mathbf{X} = \mathbf{CM}$, $\rightarrow \mathbf{H} = \mathbf{BX}$
- The fourth moments of \mathbf{H} have the form:
 $E[\mathbf{h}_i \mathbf{h}_j \mathbf{h}_k \mathbf{h}_l]$
- If the rows of \mathbf{H} were independent
 $E[\mathbf{h}_i \mathbf{h}_j \mathbf{h}_k \mathbf{h}_l] = E[\mathbf{h}_i] E[\mathbf{h}_j] E[\mathbf{h}_k] E[\mathbf{h}_l]$
- Solution: Compute \mathbf{B} such that the fourth moments of $\mathbf{H} = \mathbf{BX}$ are decoupled
 - While ensuring that \mathbf{B} is Unitary

ICA: Freeing Fourth Moments

- 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}

$$D = \begin{bmatrix} d_{11} & d_{12} & d_{13} & \dots \\ d_{21} & d_{22} & d_{23} & \dots \\ \dots & \dots & \dots & \dots \end{bmatrix}$$

- Where

$$d_{ij} = E[\sum_k \mathbf{h}_k^2 \mathbf{h}_i \mathbf{h}_j]$$

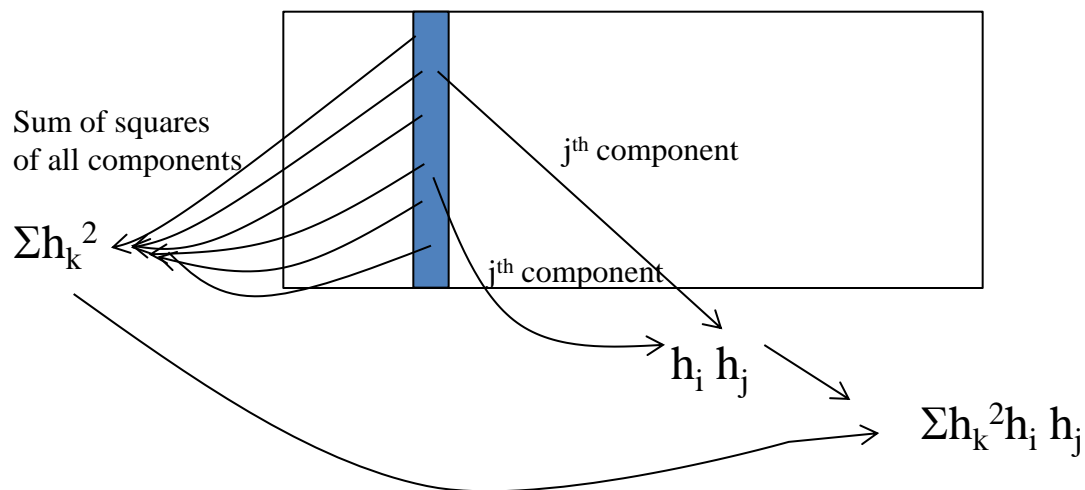
- i.e.

$$D = E[\mathbf{h}^T \mathbf{h} \mathbf{h} \mathbf{h}^T]$$

- \mathbf{h} are the columns of \mathbf{H}
- Assuming \mathbf{h} is real, else replace transposition with Hermition

ICA: The D matrix

$$D = \begin{bmatrix} d_{11} & d_{12} & d_{13} & \dots \\ d_{21} & d_{22} & d_{23} & \dots \\ \dots & \dots & \dots & \dots \end{bmatrix} \quad d_{ij} = \mathbf{E}[\sum_k \mathbf{h}_k^2 \mathbf{h}_i \mathbf{h}_j] = \frac{1}{\text{cols}(\mathbf{H})} \sum_m \sum_k h_{mk}^2 h_{mi} h_{mj}$$



- Average above term across all columns of \mathbf{H}

ICA: The D matrix

$$D = \begin{bmatrix} d_{11} & d_{12} & d_{13} & \dots \\ d_{21} & d_{22} & d_{23} & \dots \\ \dots & \dots & \dots & \dots \end{bmatrix} \quad d_{ij} = \mathbf{E}[\sum_k \mathbf{h}_k^2 \mathbf{h}_i \mathbf{h}_j] = \frac{1}{\text{cols}(\mathbf{H})} \sum_m \sum_k h_{mk}^2 h_{mi} h_{mj}$$

- If the \mathbf{h}_i terms were independent

- For $i \neq j$

$$E\left[\sum_k \mathbf{h}_k^2 \mathbf{h}_i \mathbf{h}_j\right] = E[\mathbf{h}_i^3]E[\mathbf{h}_j] + E[\mathbf{h}_j^3]E[\mathbf{h}_i] + \sum_{k \neq i, k \neq j} E[\mathbf{h}_k^2]E[\mathbf{h}_i]E[\mathbf{h}_j]$$

- Centered: $E[\mathbf{h}_j] = 0 \rightarrow E[\sum_k \mathbf{h}_k^2 \mathbf{h}_i \mathbf{h}_j] = 0$ for $i \neq j$

- For $i = j$

$$E\left[\sum_k \mathbf{h}_k^2 \mathbf{h}_i \mathbf{h}_j\right] = E[\mathbf{h}_i^4] + E[\mathbf{h}_i^2] \sum_{k \neq i} E[\mathbf{h}_k^2] \neq 0$$

- Thus, if the \mathbf{h}_i terms were independent, $d_{ij} = 0$ if $i \neq j$
- i.e., if \mathbf{h}_i were independent, D would be a diagonal matrix
 - **Let us diagonalize D**

Diagonalizing D

- Compose a fourth order matrix from \mathbf{X}
 - Recall: $\mathbf{X} = \mathbf{C}\mathbf{M}$, $\mathbf{H} = \mathbf{B}\mathbf{X} = \mathbf{B}\mathbf{C}\mathbf{M}$
 - \mathbf{B} is what we're trying to learn to make \mathbf{H} independent
- Note: if $\mathbf{H} = \mathbf{B}\mathbf{X}$, then each $\mathbf{h} = \mathbf{B}\mathbf{x}$
- The fourth moment matrix of \mathbf{H} is
- $$\begin{aligned}\mathbf{D} &= E[\mathbf{h}^T \mathbf{h} \mathbf{h} \mathbf{h}^T] = E[\mathbf{x}^T \mathbf{B}\mathbf{B}^T \mathbf{x} \mathbf{B}^T \mathbf{x} \mathbf{x}^T \mathbf{B}] \\ &= E[\mathbf{x}^T \mathbf{x} \mathbf{B}^T \mathbf{x} \mathbf{x}^T \mathbf{B}] \\ &= \mathbf{B}^T E[\mathbf{x}^T \mathbf{x} \mathbf{x} \mathbf{x}^T] \mathbf{B}\end{aligned}$$

Diagonalizing D

- Objective: Estimate \mathbf{B} such that the fourth moment of $\mathbf{H} = \mathbf{B}\mathbf{X}$ is diagonal
- Compose $\mathbf{D}_x = E[\mathbf{x}^T \mathbf{x} \mathbf{x} \mathbf{x}^T]$
- Diagonalize \mathbf{D}_x via Eigen decomposition
$$\mathbf{D}_x = \mathbf{U}\mathbf{\Lambda}\mathbf{U}^T$$
- $\mathbf{B} = \mathbf{U}^T$
 - That's it!!!!

B frees the fourth moment

$$\mathbf{D}_x = \mathbf{U}\Lambda\mathbf{U}^T ; \quad \mathbf{B} = \mathbf{U}^T$$

- \mathbf{U} is a unitary matrix, i.e. $\mathbf{U}^T\mathbf{U} = \mathbf{U}\mathbf{U}^T = \mathbf{I}$ (identity)
- $\mathbf{H} = \mathbf{B}\mathbf{X} = \mathbf{U}^T\mathbf{X}$
- $\mathbf{h} = \mathbf{U}^T\mathbf{x}$
- The fourth moment matrix of H is
$$\begin{aligned} \mathbf{E}[\mathbf{h}^T \mathbf{h} \mathbf{h} \mathbf{h}^T] &= \mathbf{U}^T \mathbf{E}[\mathbf{x}^T \mathbf{x} \mathbf{x} \mathbf{x}^T] \mathbf{U} \\ &= \mathbf{U}^T \mathbf{D}_x \mathbf{U} \\ &= \mathbf{U}^T \mathbf{U} \Lambda \mathbf{U}^T \mathbf{U} = \Lambda \end{aligned}$$
- The fourth moment matrix of $\mathbf{H} = \mathbf{U}^T\mathbf{X}$ is Diagonal!!

Overall Solution

- $\mathbf{H} = \mathbf{A}\mathbf{M} = \mathbf{B}\mathbf{C}\mathbf{M}$
 - \mathbf{C} is the (transpose of the) matrix of Eigen vectors of $\mathbf{M}\mathbf{M}^T$
- $\mathbf{X} = \mathbf{C}\mathbf{M}$
- $\mathbf{A} = \mathbf{B}\mathbf{C} = \mathbf{U}^T\mathbf{C}$
 - \mathbf{B} is the (transpose of the) matrix of Eigenvectors of $\mathbf{X} \cdot \mathit{diag}(\mathbf{X}^T\mathbf{X}) \cdot \mathbf{X}^T$

Independent Component Analysis

- Goal: to derive a matrix \mathbf{A} such that the rows of \mathbf{AM} are independent
- Procedure:
 1. “Center” \mathbf{M}
 2. Compute the autocorrelation matrix \mathbf{R}_{MM} of \mathbf{M}
 3. Compute whitening matrix \mathbf{C} via Eigen decomposition
$$\mathbf{R}_{MM} = \mathbf{E}\mathbf{S}\mathbf{E}^T, \quad \mathbf{C} = \mathbf{S}^{-1/2}\mathbf{E}^T$$
 4. Compute $\mathbf{X} = \mathbf{CM}$
 5. Compute the fourth moment matrix $\mathbf{D}' = E[\mathbf{x}^T \mathbf{x} \mathbf{x} \mathbf{x}^T]$
 6. Diagonalize \mathbf{D}' via Eigen decomposition
 7. $\mathbf{D}' = \mathbf{U}\mathbf{\Lambda}\mathbf{U}^T$
 8. Compute $\mathbf{A} = \mathbf{U}^T \mathbf{C}$
- The fourth moment matrix of $\mathbf{H} = \mathbf{AM}$ is diagonal
 - Note that the autocorrelation matrix of \mathbf{H} will also be diagonal

ICA by diagonalizing moment matrices

- The procedure just outlined, while fully functional, has shortcomings
 - Only a subset of fourth order moments are considered
 - There are many other ways of constructing fourth-order moment matrices that would ideally be diagonal
 - Diagonalizing the particular fourth-order moment matrix we have chosen is not guaranteed to diagonalize every other fourth-order moment matrix
- JADE: (Joint Approximate Diagonalization of Eigenmatrices), J.F. Cardoso
 - Jointly diagonalizes several fourth-order moment matrices
 - More effective than the procedure shown, but computationally more expensive

Enforcing Independence

- Specifically ensure that the components of \mathbf{H} are independent
 - $\mathbf{H} = \mathbf{A}\mathbf{M}$
- *Contrast function*: A non-linear function that has a minimum value when the *output components* are independent
- Define and minimize a contrast function
 - » $F(\mathbf{A}\mathbf{M})$
- Contrast functions are often only *approximations* too..

A note on pre-whitening

- The mixed signal is usually “prewhitened”
 - Normalize variance along all directions
 - Eliminate second-order dependence
- Eigen decomposition $\mathbf{M}\mathbf{M}^T = \mathbf{E}\mathbf{S}\mathbf{E}^T$
- $\mathbf{C} = \mathbf{S}^{-1/2}\mathbf{E}^T$
- Can use *first K* columns of \mathbf{E} only if only K independent sources are expected
 - In microphone array setup – only $K < M$ sources
- $\mathbf{X} = \mathbf{C}\mathbf{M}$
 - $E[\mathbf{x}_i\mathbf{x}_j] = \delta_{ij}$ for centered signal

The contrast function

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

$$I(\mathbf{H}) = \sum_i H(\bar{\mathbf{h}}_i) - H(\bar{\mathbf{h}})$$

- With constraint : $\mathbf{H} = \mathbf{B}\mathbf{X}$
 - \mathbf{X} is “whitened” \mathbf{M}

Linear Functions

- $\mathbf{h} = \mathbf{B}\mathbf{x}$, $\mathbf{x} = \mathbf{B}^{-1}\mathbf{h}$
 - Individual columns of the \mathbf{H} and \mathbf{X} matrices
 - \mathbf{x} is mixed signal, \mathbf{B} is the *unmixing* matrix

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

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

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

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

The contrast function

$$I(\mathbf{H}) = \sum_i H(\bar{\mathbf{h}}_i) - H(\bar{\mathbf{H}})$$

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

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

$$J(\mathbf{H}) = \sum_i H(\bar{\mathbf{h}}_i) - \log |\mathbf{B}|$$

- Minimize the above to obtain \mathbf{B}

An alternate approach

- Definition of Independence – if x and y are independent:
 - $E[f(x)g(y)] = E[f(x)]E[g(y)]$
 - Must hold for *every* $f()$ and $g()!!$

An alternate approach

- Define $\mathbf{g}(\mathbf{H}) = \mathbf{g}(\mathbf{BX})$ (component-wise function)

$g(h_{11})$	$g(h_{21})$...
$g(h_{12})$	$g(h_{22})$	
•	•	
•	•	
•	•	

- Define $\mathbf{f}(\mathbf{H}) = \mathbf{f}(\mathbf{BX})$

$f(h_{11})$	$f(h_{21})$...
$f(h_{12})$	$f(h_{22})$	
•	•	
•	•	
•	•	

An alternate approach

- $\mathbf{P} = \mathbf{g}(\mathbf{H}) \mathbf{f}(\mathbf{H})^T = \mathbf{g}(\mathbf{B}\mathbf{X}) \mathbf{f}(\mathbf{B}\mathbf{X})^T$

$$\mathbf{P} = \begin{array}{|ccc|} \hline P_{11} & P_{21} & \dots \\ P_{12} & P_{22} & \\ \cdot & \cdot & \\ \cdot & \cdot & \\ \cdot & \cdot & \\ \hline \end{array}$$

$$\mathbf{P}_{ij} = \mathbf{E}[\mathbf{g}(h_i)\mathbf{f}(h_j)]$$

This is a square matrix

- Must ideally be

$$\mathbf{Q} = \begin{array}{|ccc|} \hline Q_{11} & Q_{21} & \dots \\ Q_{12} & Q_{22} & \\ \cdot & \cdot & \\ \cdot & \cdot & \\ \cdot & \cdot & \\ \hline \end{array}$$

$$Q_{ij} = E[\mathbf{g}(h_i)]E[\mathbf{f}(h_j)] \quad i \neq j$$

$$Q_{ii} = E[\mathbf{g}(h_i)\mathbf{f}(h_i)]$$

- Error = $\|\mathbf{P}-\mathbf{Q}\|_F^2$

An alternate approach

- Ideal value for \mathbf{Q}

$$\mathbf{Q} = \begin{array}{|ccc|} \hline Q_{11} & Q_{21} & \dots \\ \hline Q_{12} & Q_{22} & \\ \hline \cdot & \cdot & \\ \hline \cdot & \cdot & \\ \hline \cdot & \cdot & \\ \hline \end{array} \quad \begin{array}{l} Q_{ij} = E[g(h_i)]E[f(h_j)] \quad i \neq j \\ \\ Q_{ii} = E[g(h_i)f(h_i)] \end{array}$$

- If $g()$ and $h()$ are odd symmetric functions
 $E[g(h_i)] = 0$ for all i
 - Since $E[h_i] = 0$ (\mathbf{H} is centered)
 - \mathbf{Q} is a Diagonal Matrix!!!

An alternate approach

- Minimize Error

$$\mathbf{P} = \mathbf{g}(\mathbf{B}\mathbf{X})\mathbf{f}(\mathbf{B}\mathbf{X})^T$$

$$\mathbf{Q} = \textit{Diagonal}$$

$$\textit{error} = \|\mathbf{P} - \mathbf{Q}\|_F^2$$

- Leads to trivial Widrow Hopf type iterative rule:

$$\mathbf{E} = \textit{Diag} - \mathbf{g}(\mathbf{B}\mathbf{X})\mathbf{f}(\mathbf{B}\mathbf{X})^T$$

$$\mathbf{B} = \mathbf{B} + \eta\mathbf{E}\mathbf{B}^T$$

Update Rules

- Multiple solutions under different assumptions for $g()$ and $f()$
- $\mathbf{H} = \mathbf{B}\mathbf{X}$
- $\mathbf{B} = \mathbf{B} + \eta \Delta\mathbf{B}$
- Jutten Herraut : Online update
 - $\Delta B_{ij} = f(\mathbf{h}_i)g(\mathbf{h}_j)$; -- actually assumed a recursive neural network
- Bell Sejnowski
 - $\Delta\mathbf{B} = ([\mathbf{B}^T]^{-1} - \mathbf{g}(\mathbf{H})\mathbf{X}^T)$

Update Rules

- Multiple solutions under different assumptions for $g()$ and $f()$
- $\mathbf{H} = \mathbf{B}\mathbf{X}$
- $\mathbf{B} = \mathbf{B} + \eta \Delta\mathbf{B}$
- Natural gradient -- $f() = \text{identity function}$
 - $\Delta\mathbf{B} = (\mathbf{I} - g(\mathbf{H})\mathbf{H}^T)\mathbf{W}$
- Cichoki-Unbehauen
 - $\Delta\mathbf{B} = (\mathbf{I} - g(\mathbf{H})f(\mathbf{H})^T)\mathbf{W}$

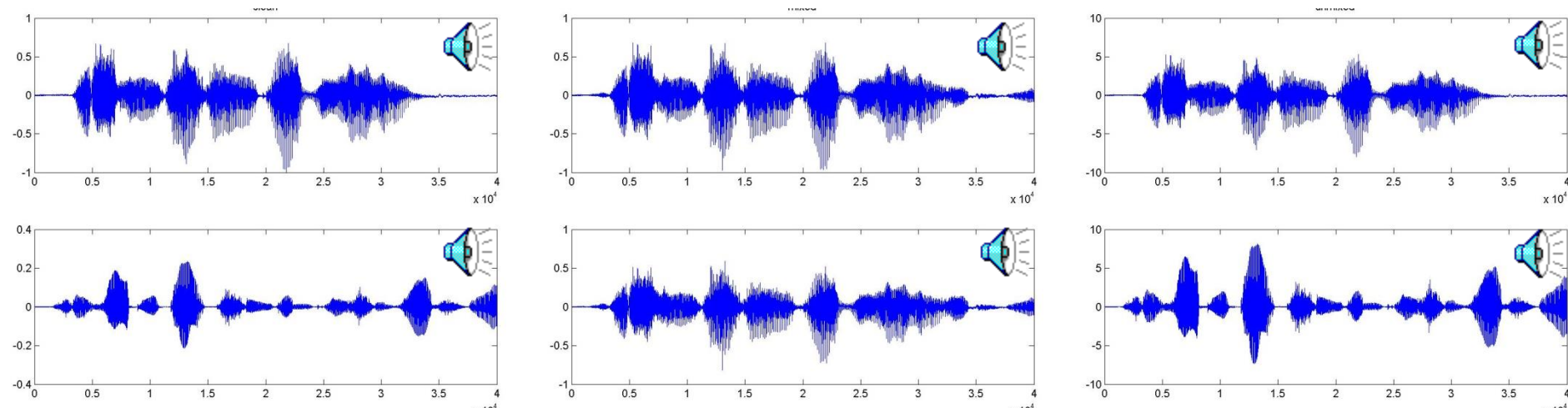
What are G() and H()

- Must be odd symmetric functions
- Multiple functions proposed

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

- Audio signals in general
 - $\Delta \mathbf{B} = (\mathbf{I} - \mathbf{H}\mathbf{H}^T - \mathbf{K}\tanh(\mathbf{H})\mathbf{H}^T)\mathbf{W}$
- Or simply
 - $\Delta \mathbf{B} = (\mathbf{I} - \mathbf{K}\tanh(\mathbf{H})\mathbf{H}^T)\mathbf{W}$

So how does it work?



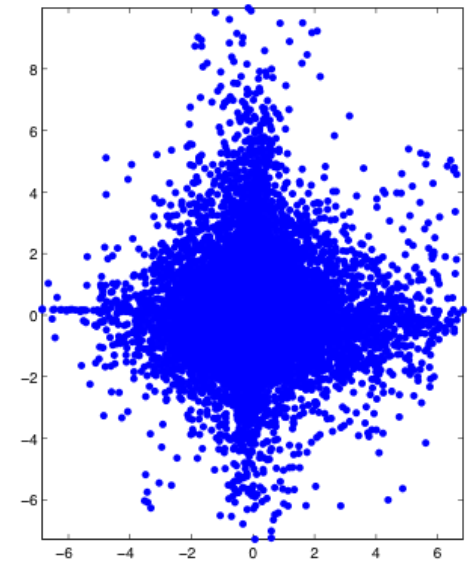
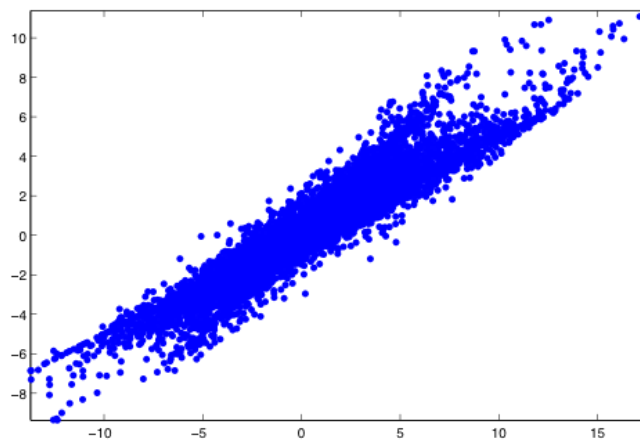
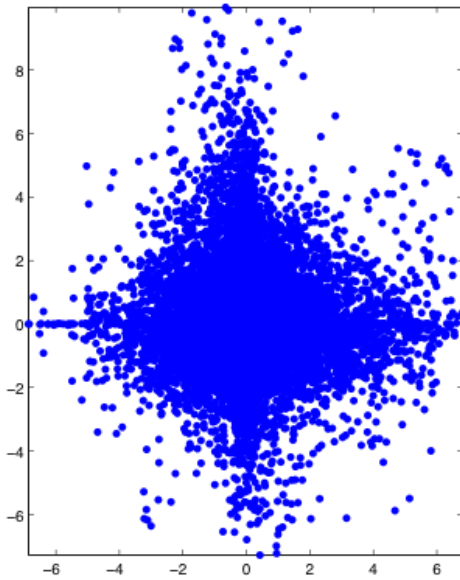
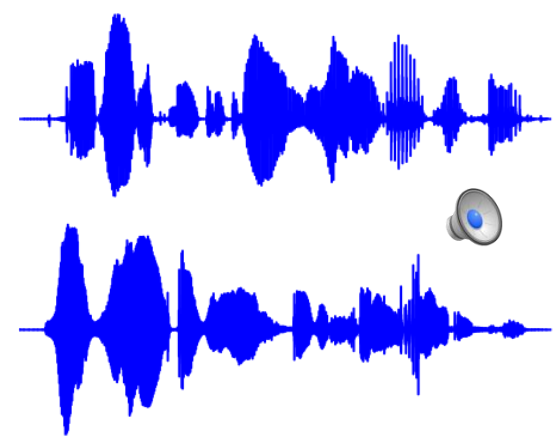
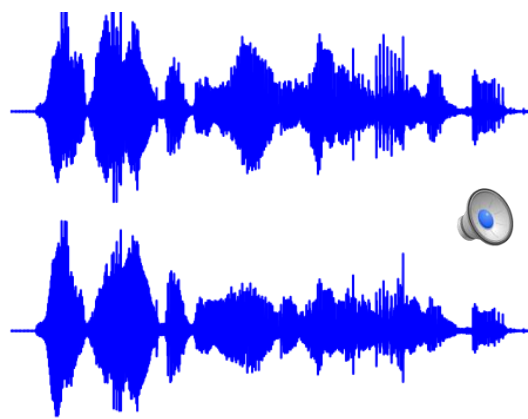
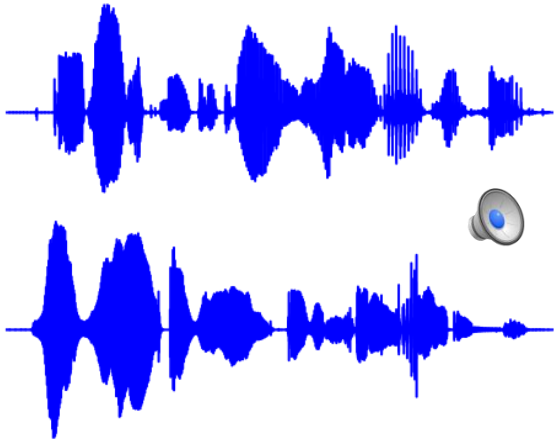
- Example with instantaneous mixture of two speakers
- Natural gradient update
- Works very well!

Another example!

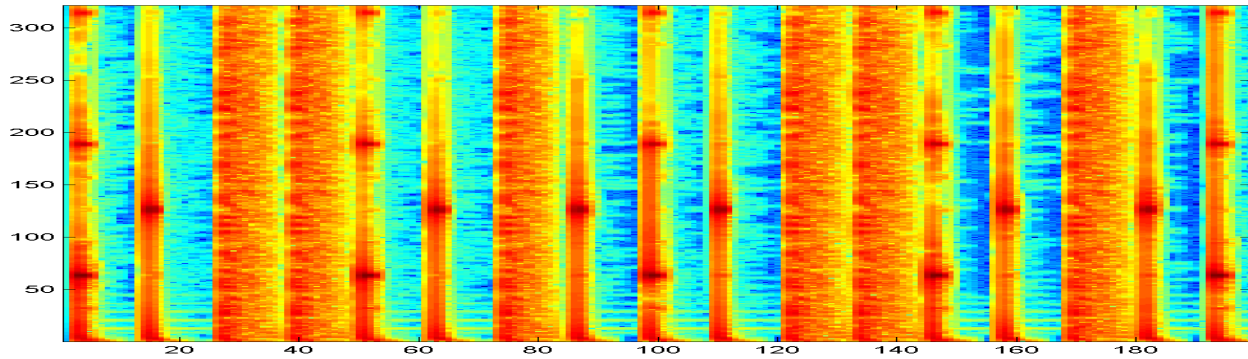
Input

Mix

Output

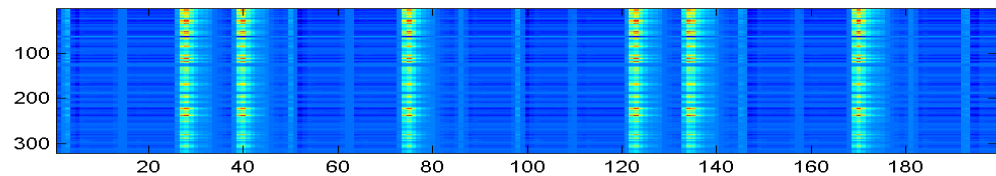
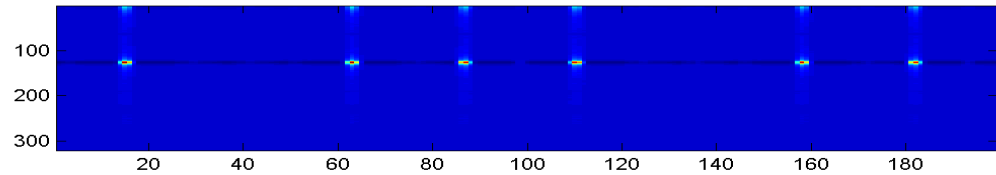
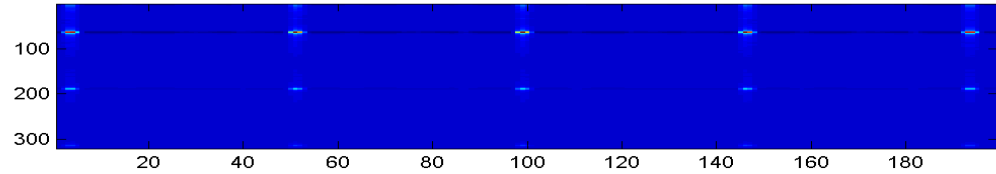
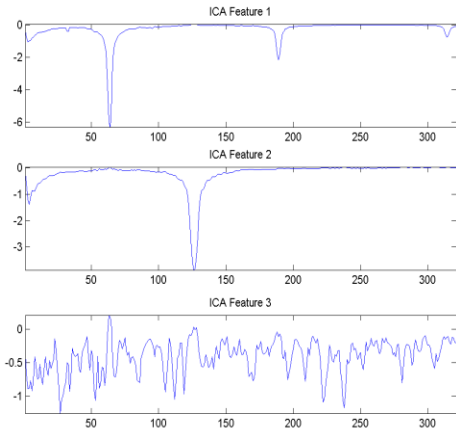


Another Example



- Three instruments..

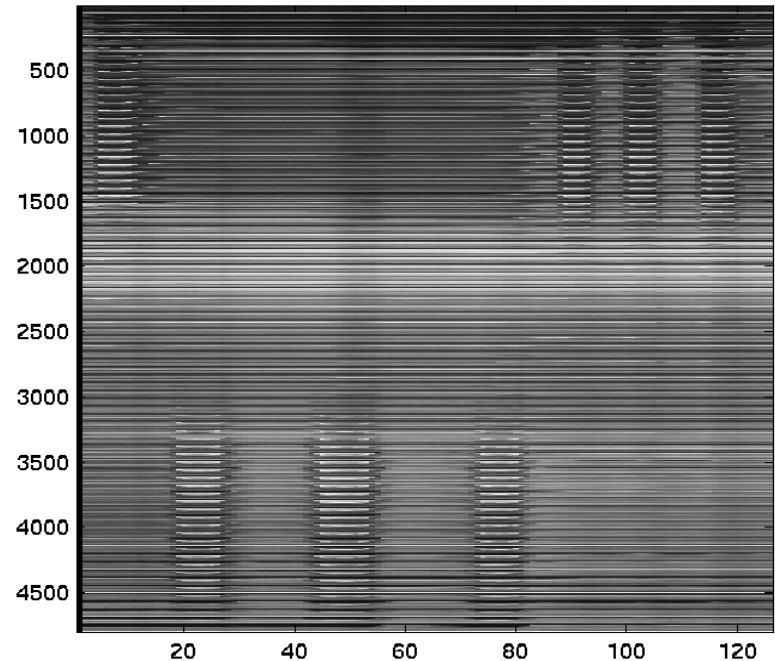
The Notes



- Three instruments..

ICA for data exploration

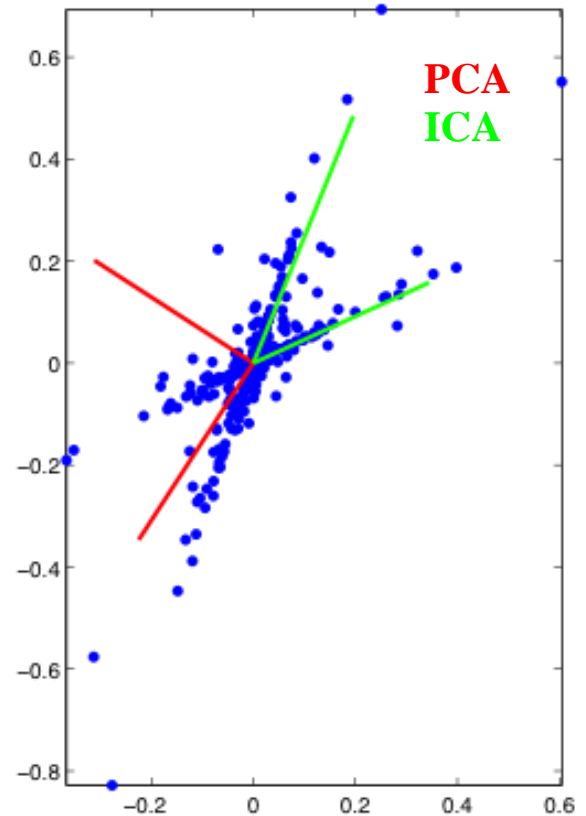
- The “bases” in PCA represent the “building blocks”
 - Ideally notes
- Very successfully used
- So can ICA be used to do the same?



ICA vs PCA bases

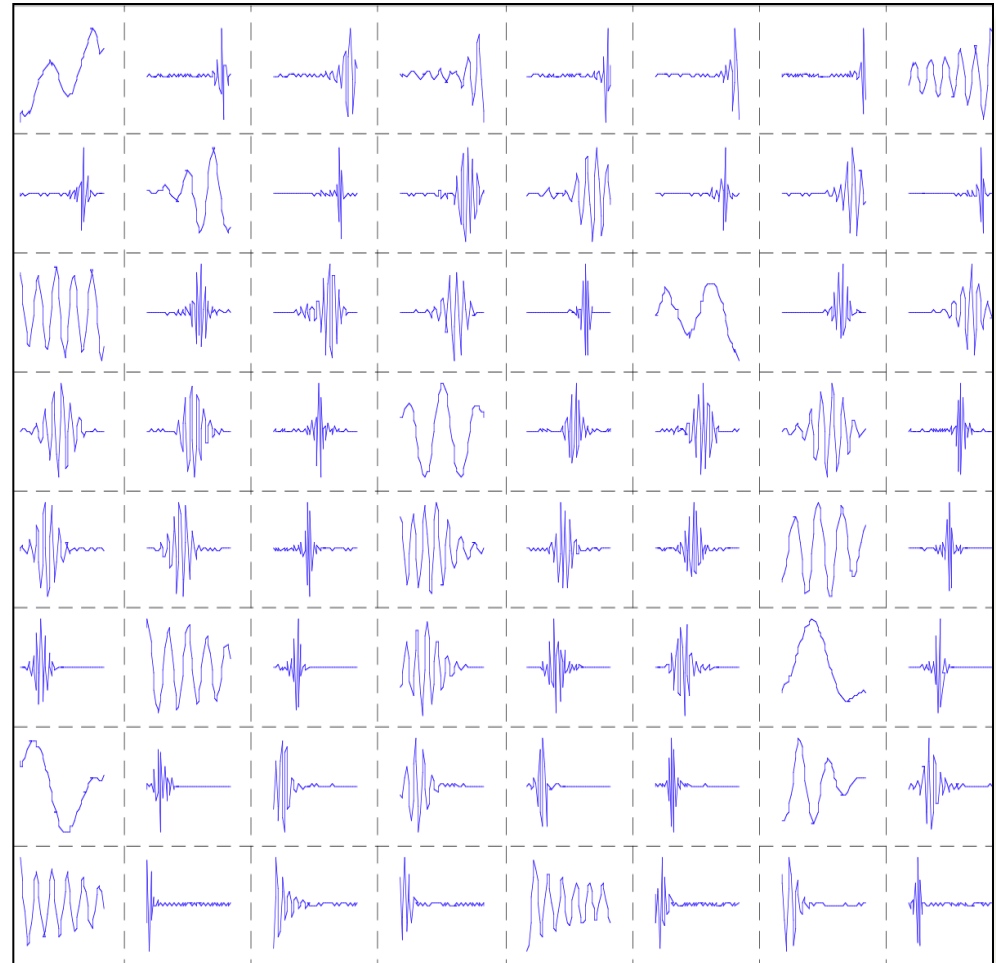
- Motivation for using ICA vs PCA
- PCA will indicate orthogonal directions of maximal variance
 - May not align with the data!
- ICA finds directions that are independent
 - More likely to “align” with the data

Non-Gaussian data



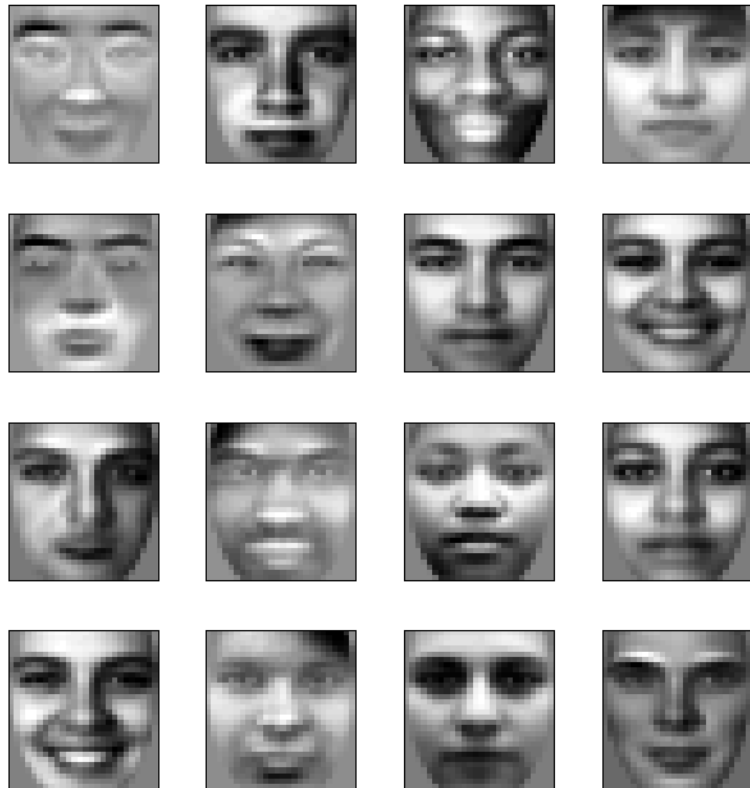
Finding useful transforms with ICA

- 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/freq localized sinusoids which is a better way to analyze sounds
- Ditto for images
 - ICA returns localizes edge filters

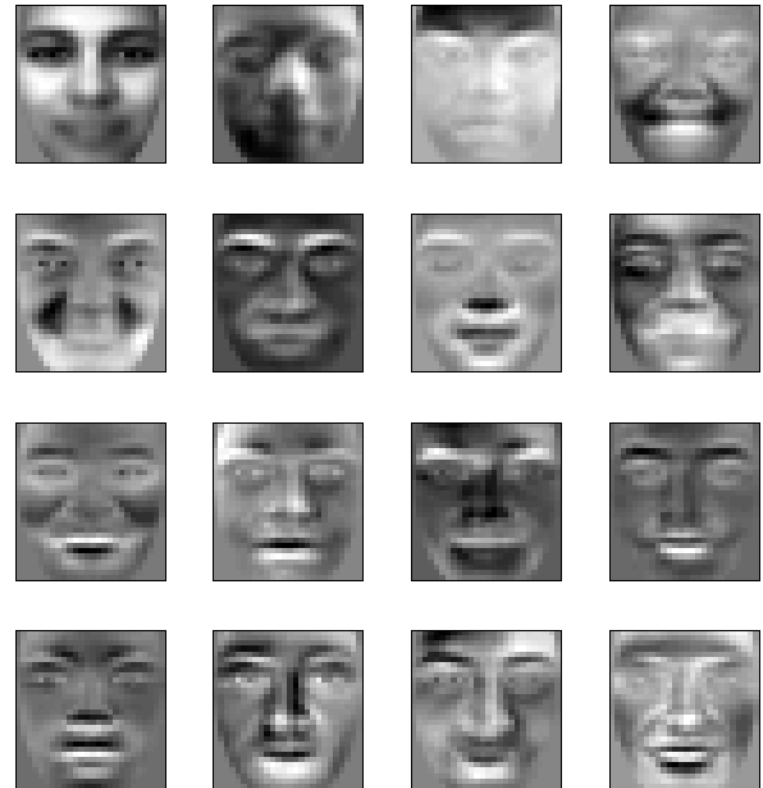


Example case: ICA-faces vs. Eigenfaces

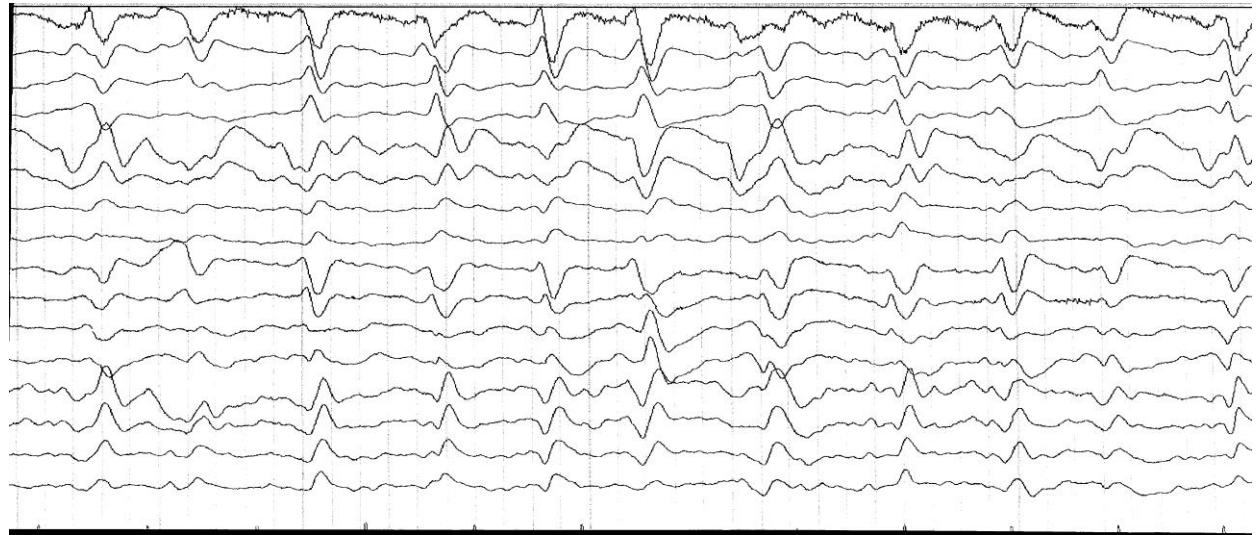
ICA-faces



Eigenfaces

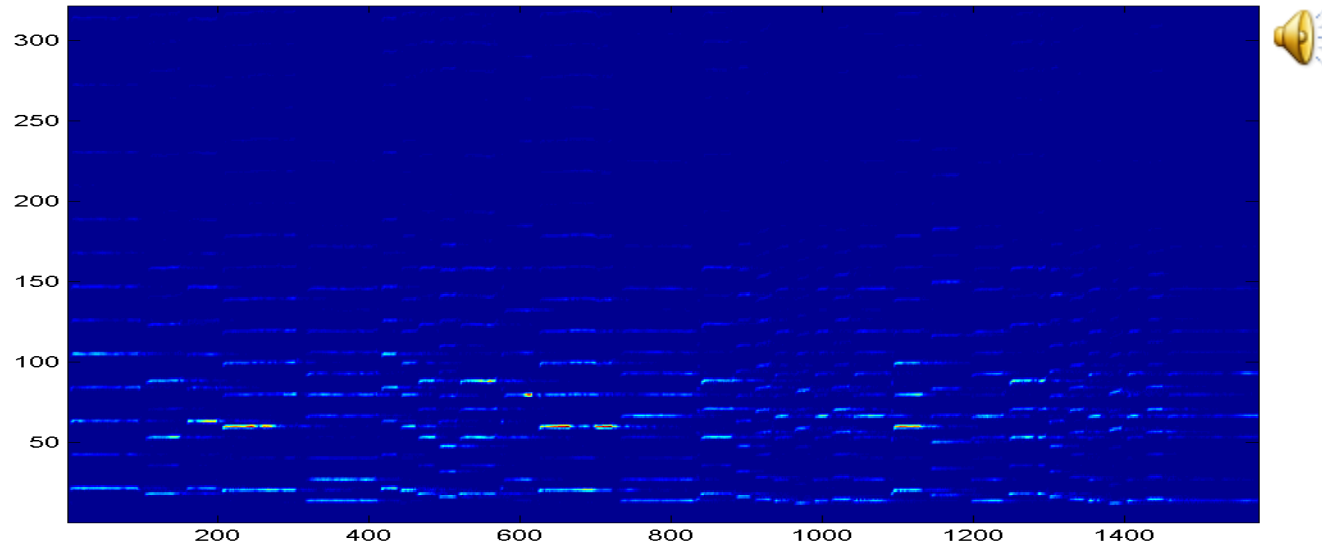


ICA for Signal Enhancement



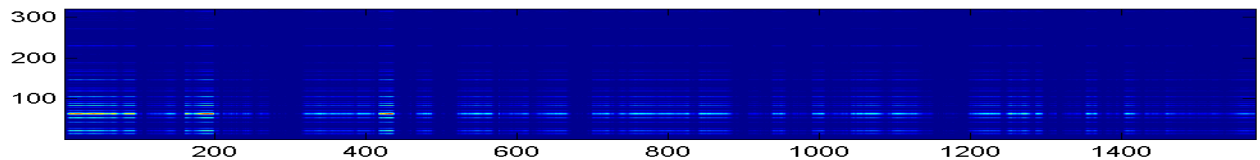
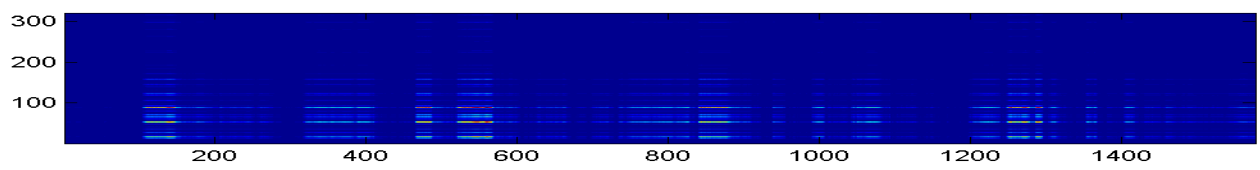
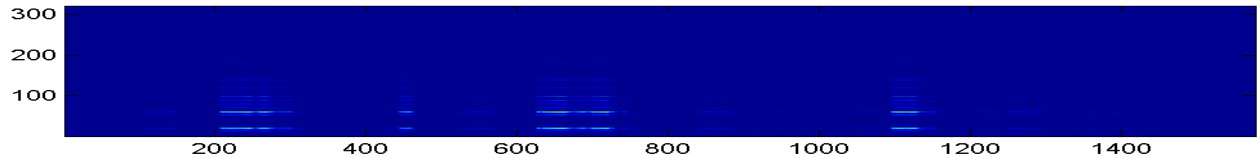
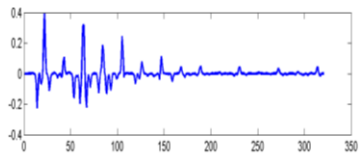
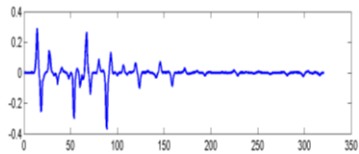
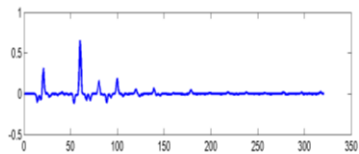
- Very commonly used to enhance EEG signals
- EEG signals are frequently corrupted by heartbeats and biorhythm signals
- ICA can be used to separate them out

So how does that work?



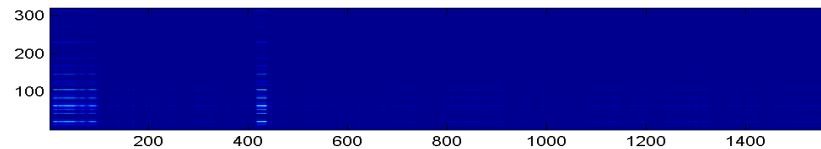
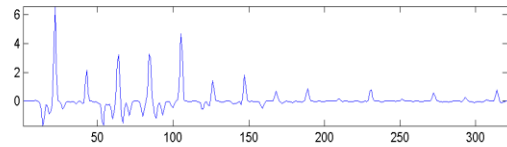
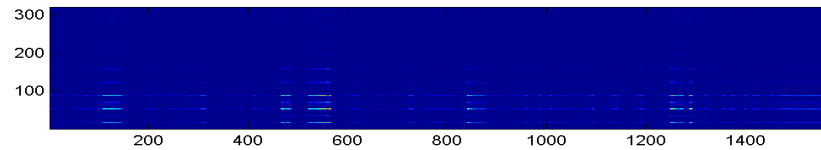
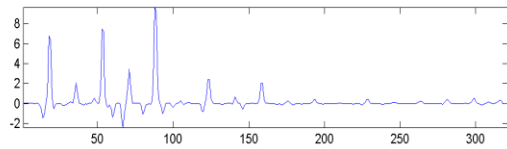
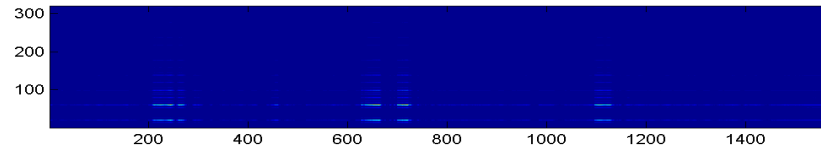
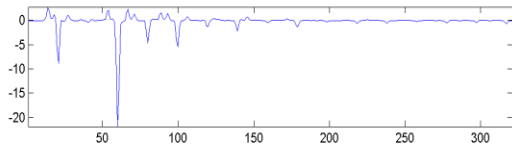
- There are 12 notes in the segment, hence we try to estimate 12 notes..

PCA solution



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

So how does this work: ICA solution



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

ICA Issues

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

What else went wrong?

- *Notes are not independent*
 - Only one note plays at a time
 - If one note plays, other notes are *not* playing
- Will deal with these later in the course..