# Machine Learning for Signal Processing
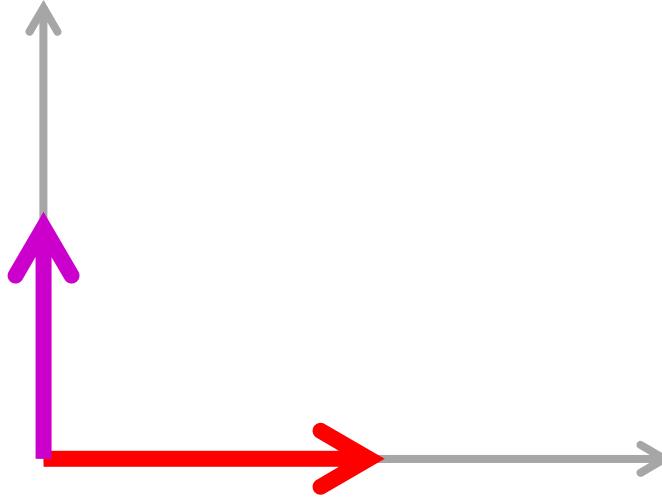# Independent Component Analysis

Instructor: Bhiksha Raj

# Revisiting the Covariance Matrix
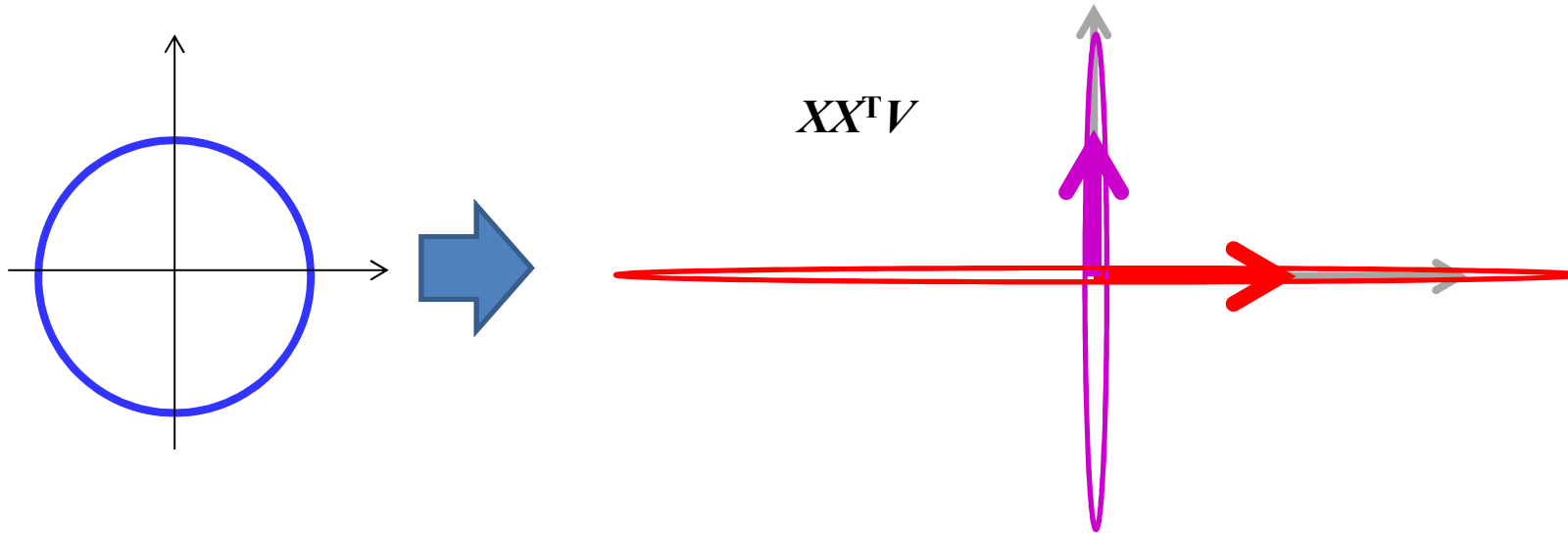
- **Assuming centered data**

- $C = \sum_X XX^\top$
- $= X_1 X_1^\top + X_2 X_2^\top + \ldots$

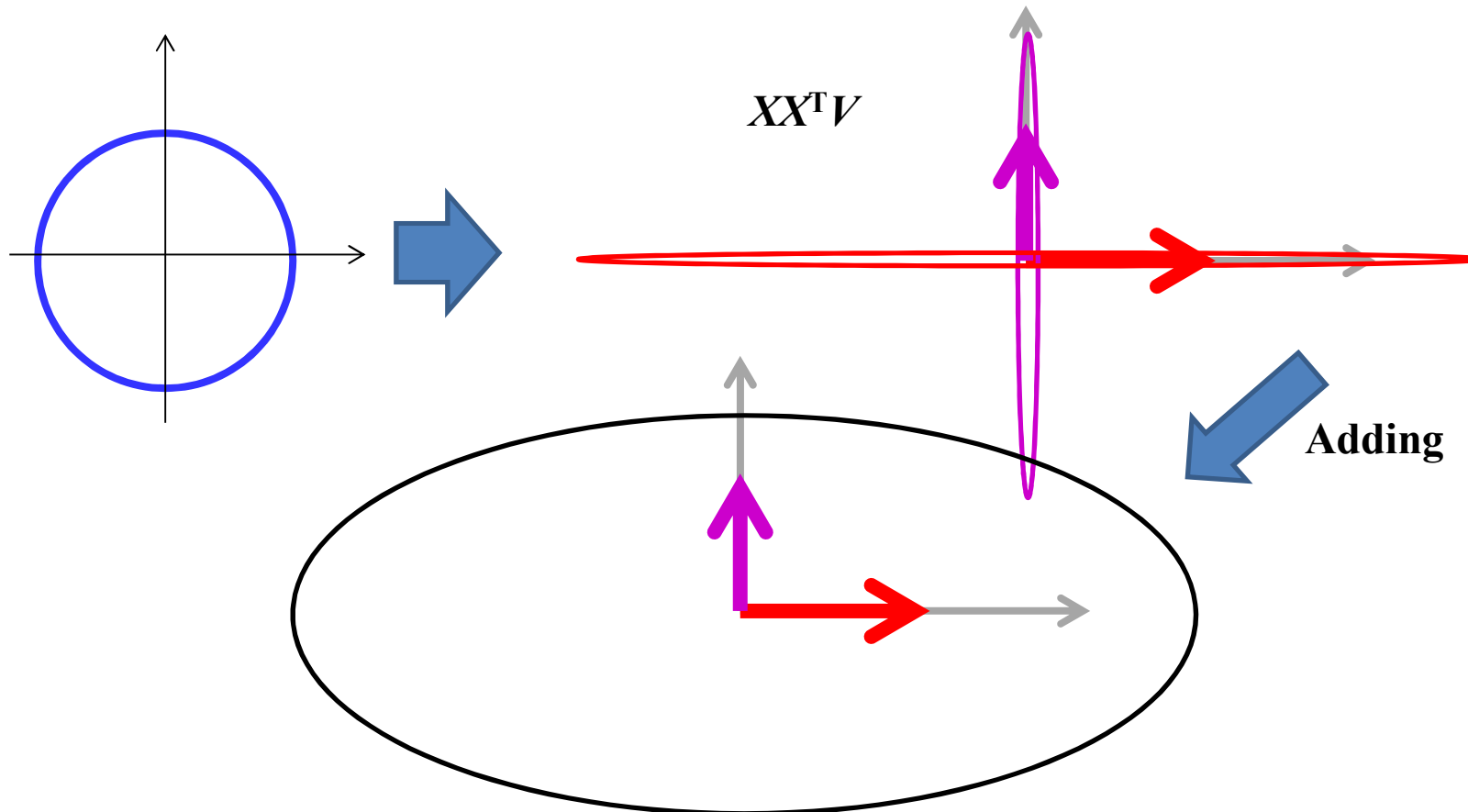- Let us view C as a transform..

# Covariance matrix as a transform



- $(X_1 X_1^{\top} + X_2 X_2^{\top} + \dots)\, V = X_1 X_1^{\top} V + X_2 X_2^{\top} V + \dots$
- Consider a 2-vector example
  - In two dimensions for illustration

# Covariance Matrix as a transform



$$XX^\mathrm{T}V$$

- Data comprises only 2 vectors..
- *Major axis of component ellipses proportional to the squared length of the corresponding vector*

# Covariance Matrix as a transform
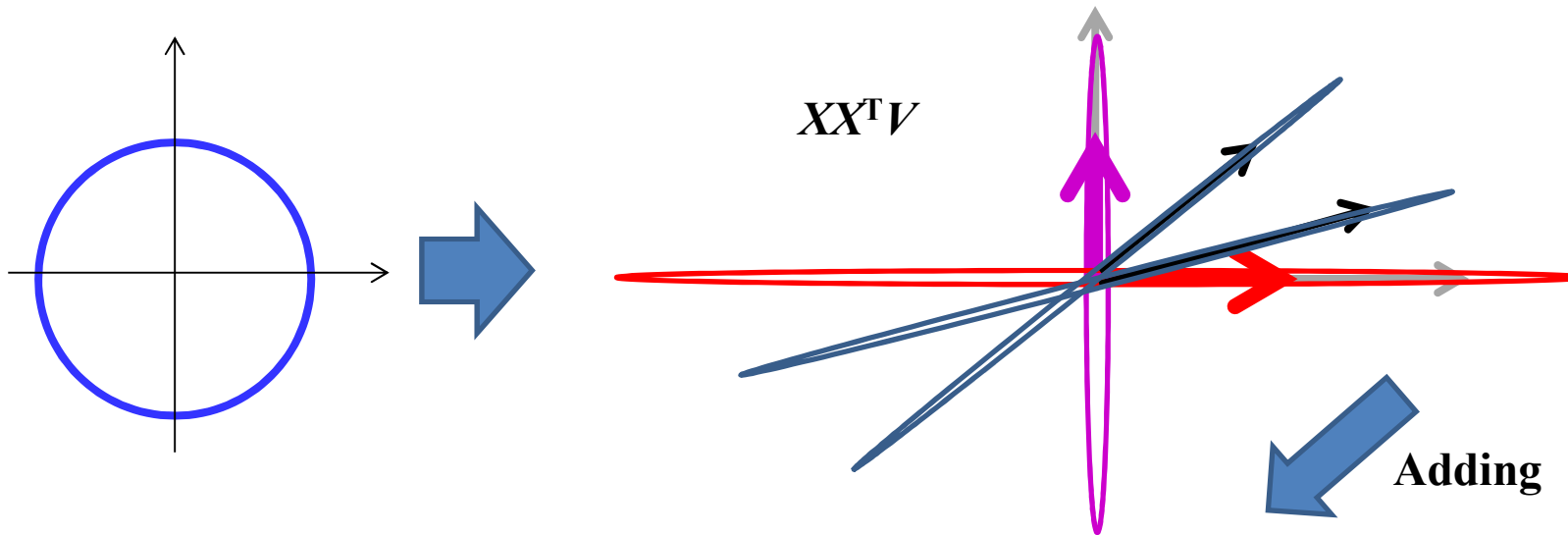


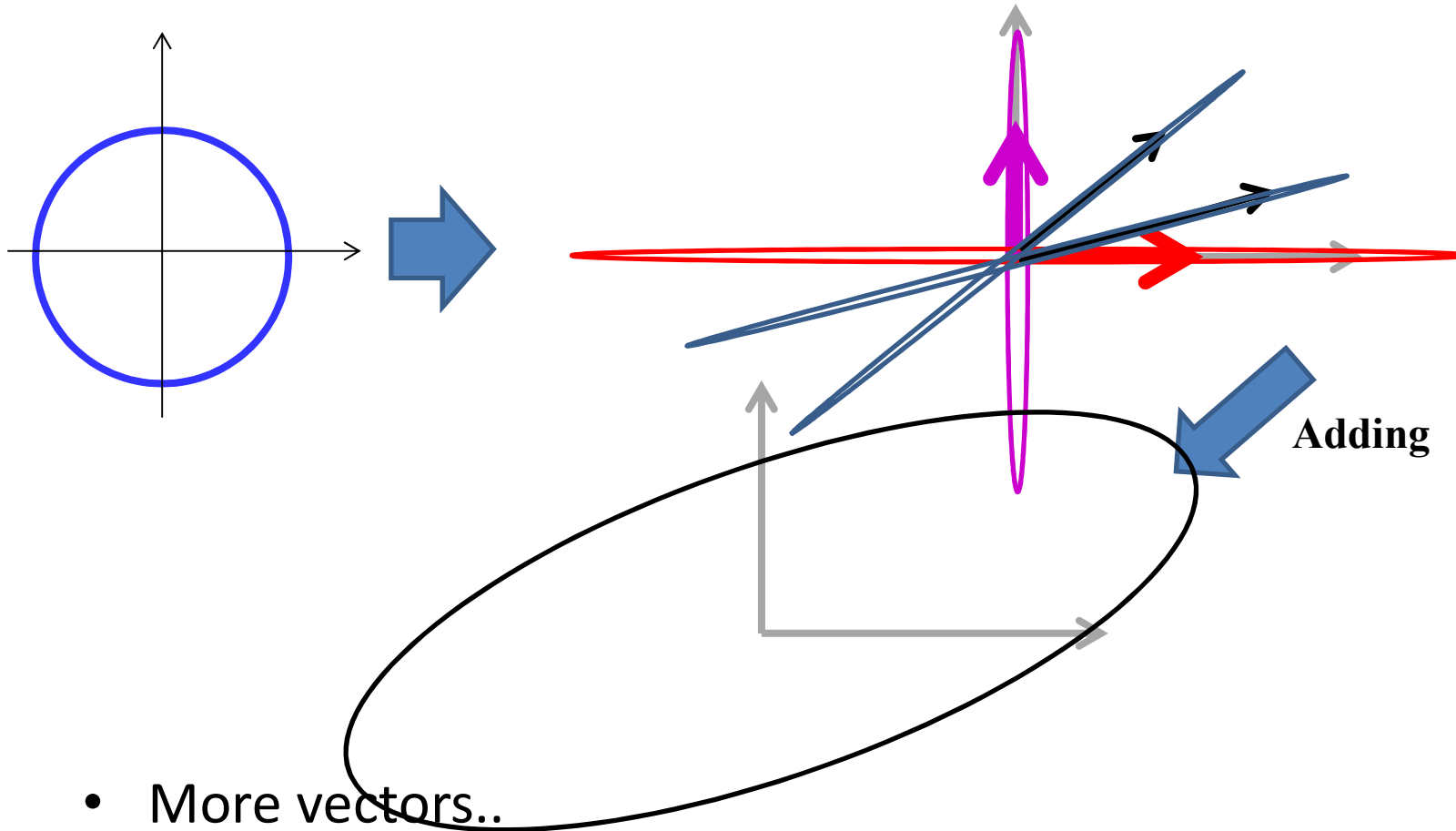$XX^\mathrm{T}V$

Adding

- Data comprises only 2 vectors..
- *Major axis of component ellipses proportional to the squared length of the corresponding vector*

# Covariance Matrix as a transform

$XX^T V$

Adding

- More vectors..
- *Major axis of component ellipses proportional to the squared length of the corresponding vector*

# Covariance Matrix as a transform



**Adding**

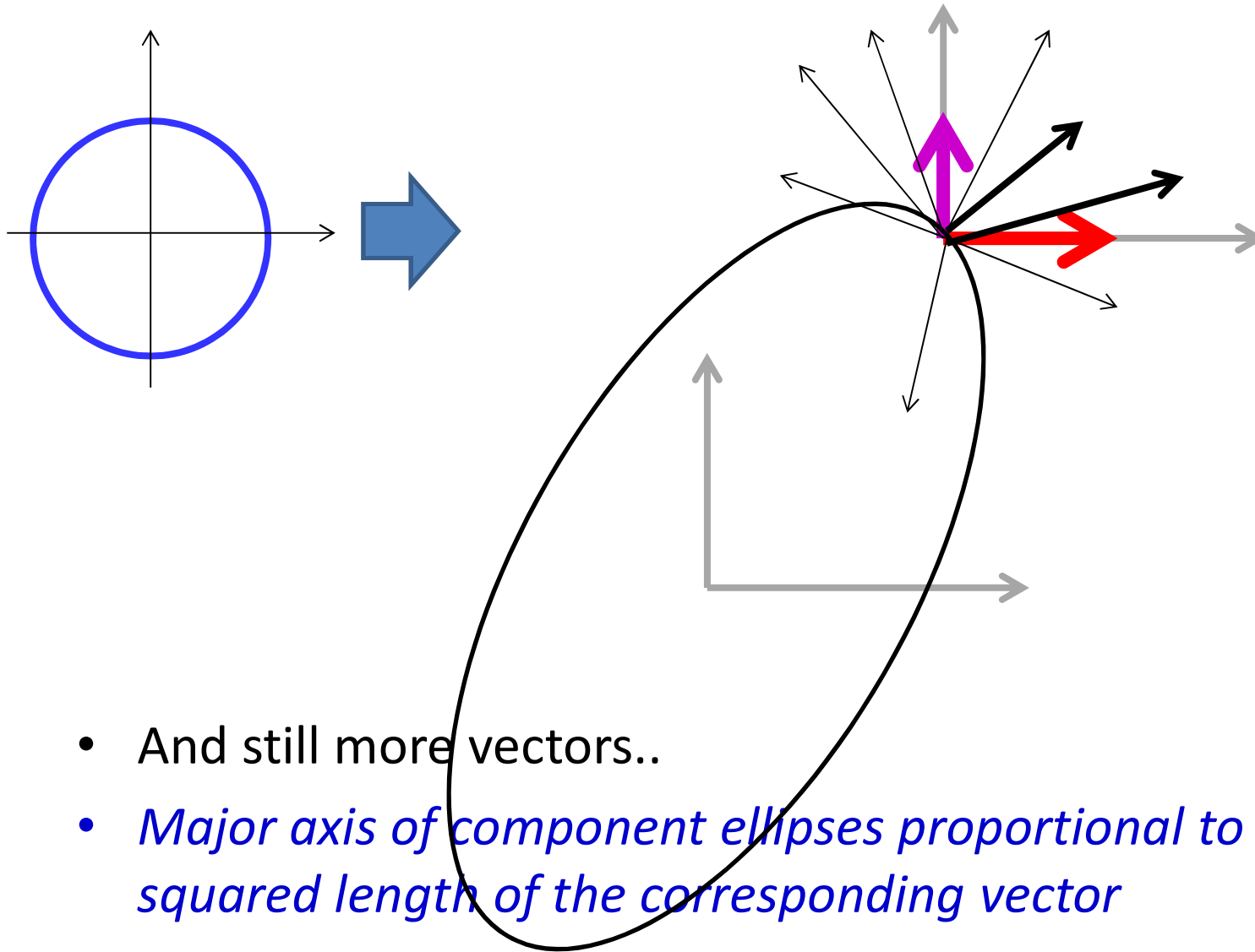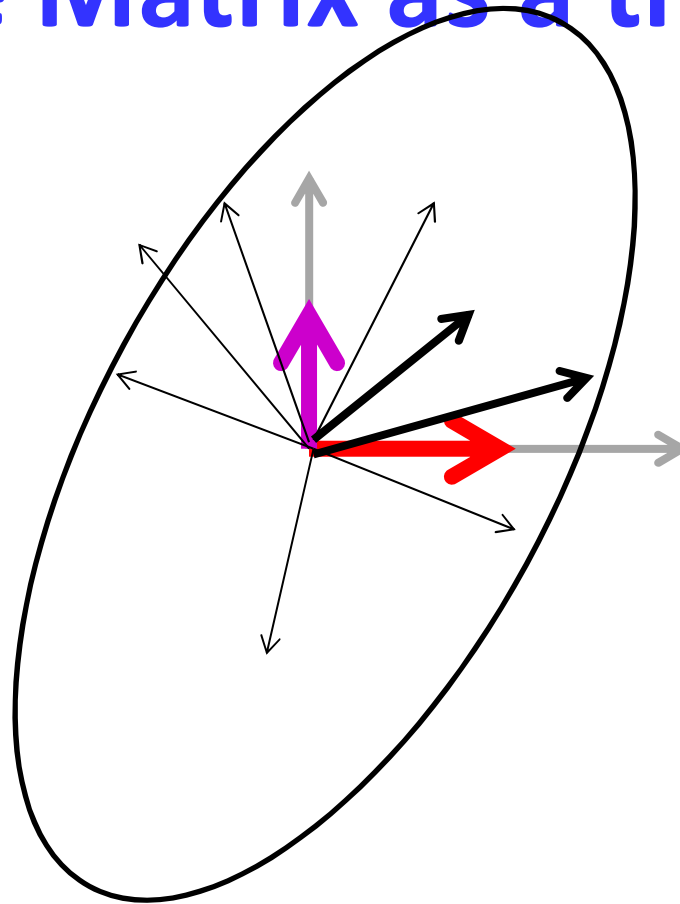- More vectors..
- *Major axis of component ellipses proportional to the squared length of the corresponding vector*

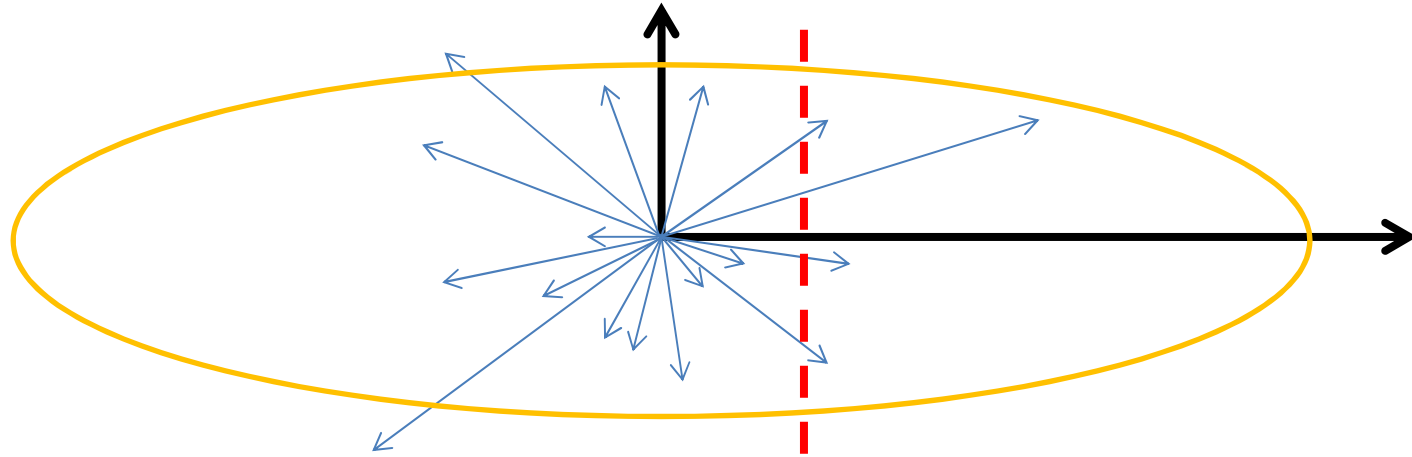# Covariance Matrix as a transform

- And still more vectors..
- *Major axis of component ellipses proportional to the squared length of the corresponding vector*

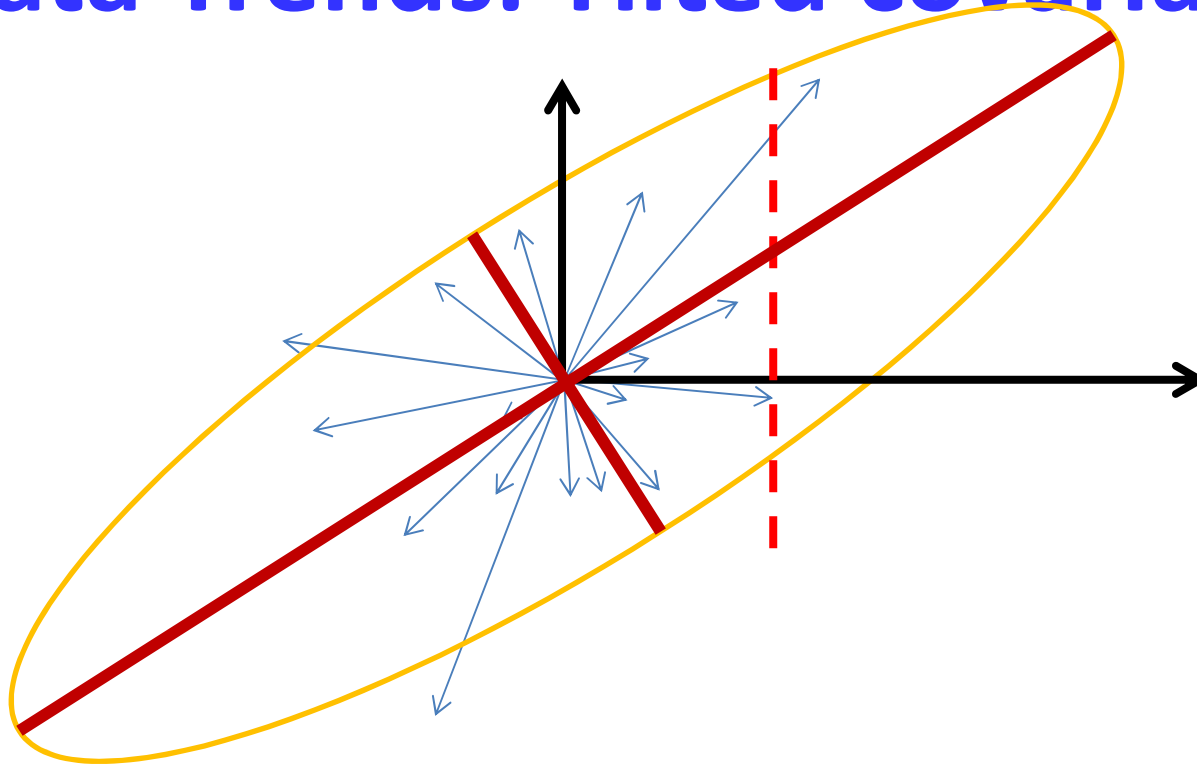# Covariance Matrix as a transform



- The covariance matrix captures the directions of maximum variance

- What does it tell us about trends?

# Data Trends: Axis aligned covariance



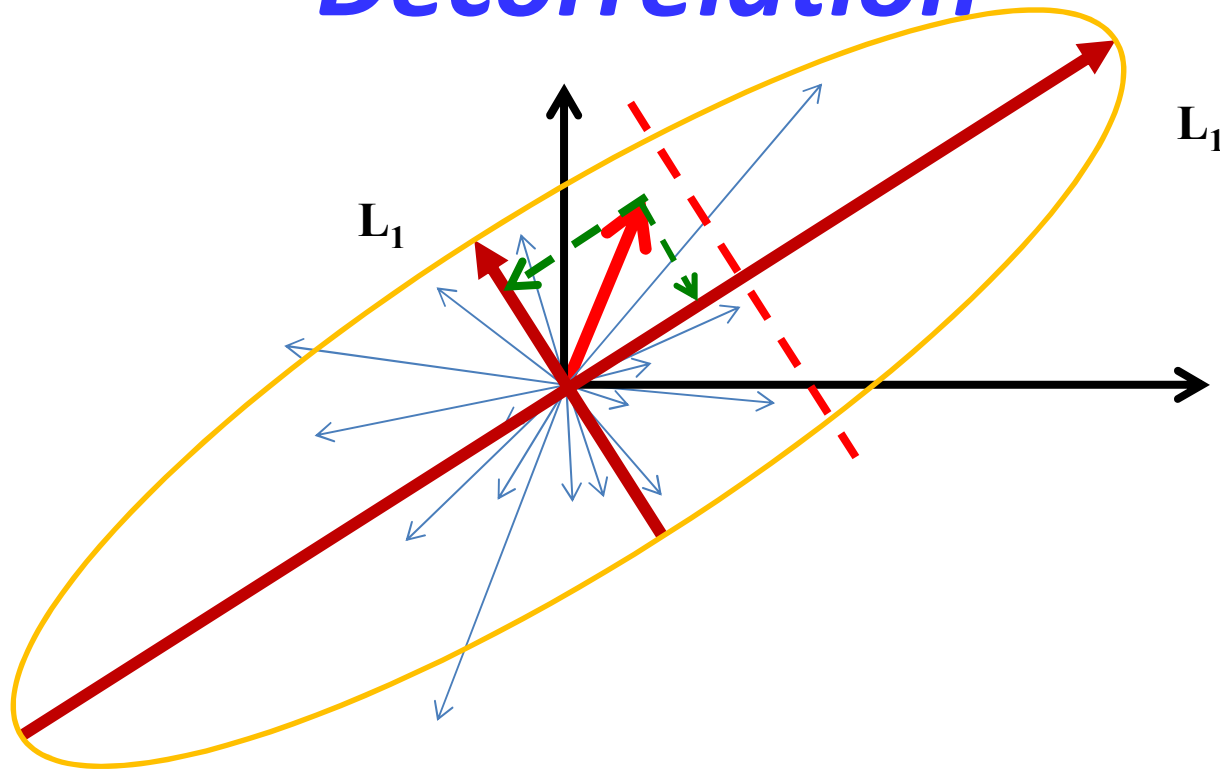- Axis aligned covariance
- At any X value, the average Y value of vectors is 0
  - X cannot predict Y
- At any Y, the average X of vectors is 0
  - Y cannot predict X
- The X and Y components are *uncorrelated*

# Data Trends: Tilted covariance



- Tilted covariance
- The average Y value of vectors at any X varies with X
  - X predicts Y
- Average X varies with Y
- The X and Y components are *correlated*

# *Decorrelation*



- Shifting to using the major axes as the coordinate system
  - $L_1$ does not predict $L_2$ and vice versa
  - In this coordinate system the data are uncorrelated
- We have ***decorrelated*** the data by rotating the axes

# The statistical concept of correlatedness

- Two variables X and Y are correlated if 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?

# Correlation vs. Causation

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

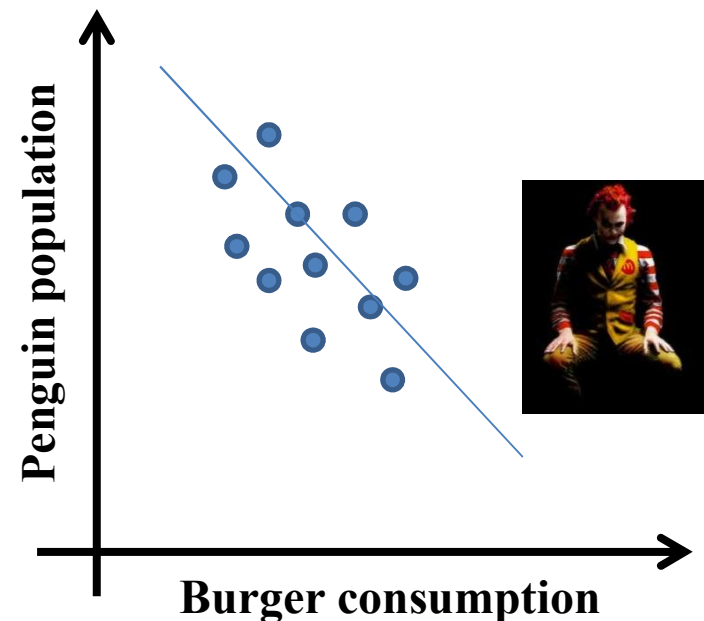- In the same period, the penguin population of Antarctica has gone down
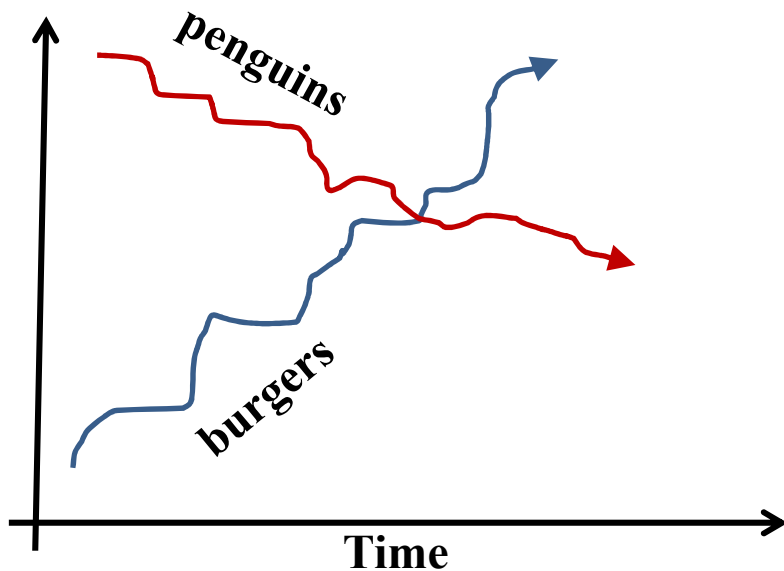
Correlation, not Causation
(unless McDonalds has a
top-secret Antarctica division)
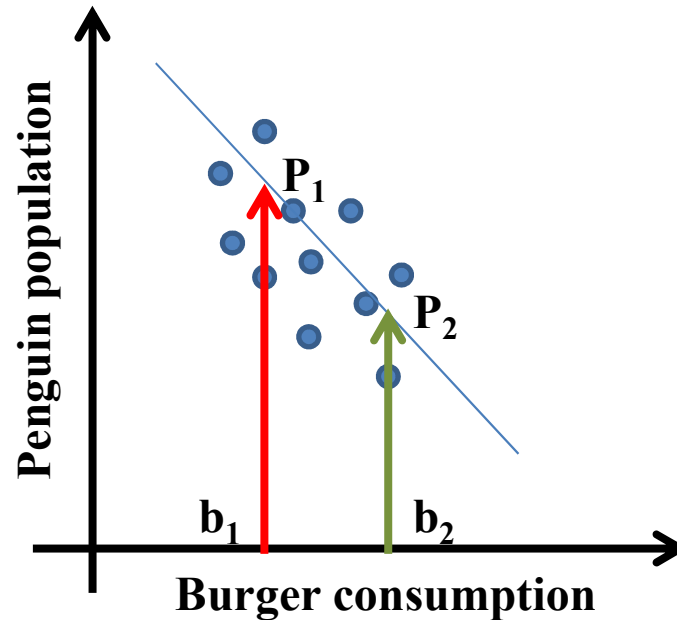
# The concept of *correlation*

- Two variables are correlated if knowing the value of one gives you information about the **expected value** of the other
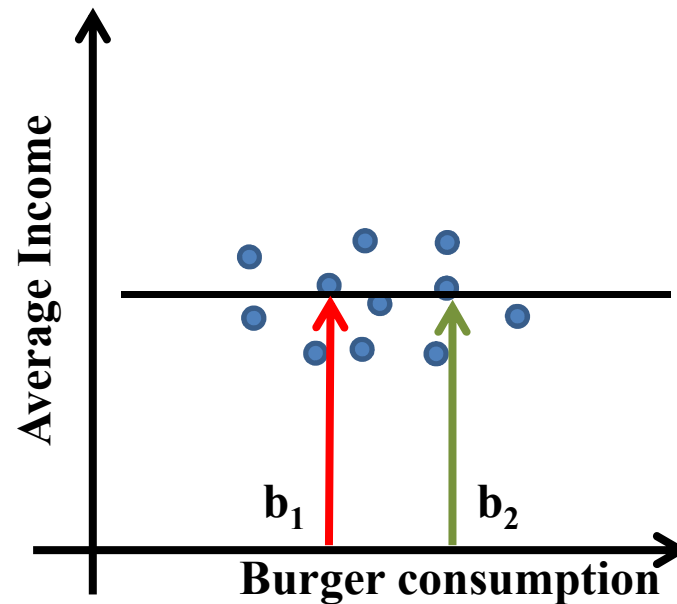
# A brief review of basic probability

- *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 in Random Variables



- Which of the above represent uncorrelated RVs?

# Benefits of uncorrelatedness..

- Uncorrelatedness of variables is generally considered desirable for modelling and analyses
  - For Euclidean error based regression models and probabilistic models, uncorrelated variables can be separately handled
    - Since the value of one doesn't affect the average value of others
    - Greatly reduces the number of model parameters
  - Otherwise their interactions must be considered

- We will frequently transform correlated variables to make them uncorrelated
  - "Decorrelating" variables

# The notion of *decorrelation*

$$\begin{pmatrix} X' \\ Y' \end{pmatrix} = M \begin{pmatrix} X \\ Y \end{pmatrix}$$

**?**

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

# What does "uncorrelated" mean

**Assuming 0 mean**



- $\mathrm{E}[X'] = $ constant
- $\mathrm{E}[Y'] = $ constant
- $\mathrm{E}[Y'|X'] = $ constant
- $\mathrm{E}[X'Y'] = \mathrm{E}[X']\,\mathrm{E}[Y']$
  - All will be 0 for centered data

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

- If $\mathbf{Y}$ is a matrix of vectors, $\mathbf{Y}\mathbf{Y}^{\mathrm{T}} = $ 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 $\mathbf{Y} = \mathbf{MX}$ such that the covariance of $\mathbf{Y}$ is diagonal
  - $\mathbf{YY}^\mathrm{T}$ is the covariance if $\mathbf{Y}$ is zero mean
  - For uncorrelated components, $\mathbf{YY}^\mathrm{T} = \mathbf{Diagonal}$

  $\Rightarrow \mathbf{MXX}^\mathrm{T}\mathbf{M}^\mathrm{T} = \mathbf{Diagonal}$

  $\Rightarrow \mathbf{M}.\mathrm{Cov}(\mathbf{X}).\mathbf{M}^\mathrm{T} = \mathbf{Diagonal}$

# Decorrelation

- Easy solution:
  - Eigen decomposition of $\mathrm{Cov}(\mathbf{X})$:
    $$\mathrm{Cov}(\mathbf{X}) = \mathbf{E}\Lambda\mathbf{E}^{\mathrm{T}}$$
  - $\mathbf{E}\mathbf{E}^{\mathrm{T}} = \mathrm{I}$
- Let $\mathbf{M} = \mathbf{E}^{\mathrm{T}}$

- $\mathbf{M}\mathrm{Cov}(\mathbf{X})\mathbf{M}^{\mathrm{T}} = \mathbf{E}^{\mathrm{T}}\mathbf{E}\Lambda\mathbf{E}^{\mathrm{T}}\mathbf{E} = \Lambda = \text{diagonal}$

- PCA: $\mathbf{Y} = \mathbf{E}^{\mathrm{T}}\mathbf{X}$
  - Projects the data onto the Eigen vectors of the covariance matrix
  - *Diagonalizes* the covariance matrix
  - "Decorrelates" the data

# PCA

$$\mathbf{X} = w_1\mathbf{E}_1 + w_2\mathbf{E}_2$$



- PCA: $\mathbf{Y} = \mathbf{E}^{\mathrm{T}}\mathbf{X}$
  - Projects the data onto the Eigen vectors of the covariance matrix
    - Changes the coordinate system to the Eigen vectors of the covariance matrix
  - *Diagonalizes* the covariance matrix
  - "Decorrelates" the data

# Decorrelating the data



- Are there other decorrelating axes?

# Decorrelating the data



- Are there other decorrelating axes?

# Decorrelating the data



- Are there other decorrelating axes?

# Decorrelating the data



- Are there other decorrelating axes?
- What about if we don't require them to be orthogonal?

# Decorrelating the data



- Are there other decorrelating axes?
- What about if we don't require them to be orthogonal?
- What is special about these axes?

# The statistical concept of *Independence*

- Two variables X and Y are *dependent* if 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 uncorrelatedness does not imply independence

# 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)

p(x)

- The expected value of an odd function of an RV is 0 if

  – The RV is 0 mean

  – The PDF is of the RV is symmetric around 0

- **E[f(X)] = 0 if f(X) is odd symmetric**

# A note on bits..

- You flip a coin.  You must inform your friend in the next room about whether the outcome was heads or tails



Digital channel

- How many bits will you have to send?

# A note on bits..

- You roll a four-side dice.  You must inform your friend in the next room about the outcome



Digital channel

- How many bits will you have to send?

# A note on bits..

- You roll an *eight-sided polyheldral* dice.  You must inform your friend in the next room about the outcome



Digital channel

- How many bits will you have to send?

# A note on bits..

- You roll a *six-sided* dice.  You must inform your friend in the next room about the outcome



Digital channel

- How many bits will you have to send?

# Batching up 6-sided dice rolls



- Instead of sending individual rolls, you roll the dice *twice*
  - And send the *pair* to your friend
- How many bits do you send *per roll?*

| 1 | 1 |
|---|---|
| 1 | 2 |
| 1 | 3 |
| .. | .. |
| 2 | 1 |
| 2 | 2 |
| .. | .. |
| 6 | 6 |

# Batching up 6-sided dice rolls



- Instead of sending individual rolls, you roll the dice *twice*
  - And send the *pair* to your friend
- How many bits do you send *per roll?*
- 36 combinations: 6 bits per pair of numbers
  - Still 3 bits per roll

| | |
|---|---|
| 1 | 1 |
| 1 | 2 |
| 1 | 3 |
| .. | .. |
| 2 | 1 |
| 2 | 2 |
| .. | .. |
| 6 | 6 |

# Batching up 6-sided dice rolls



Digital channel

| 1 | 1 | 1 |
|---|---|---|
| 1 | 1 | 2 |
| .. | .. | .. |
| 1 | 6 | 3 |
| .. |  | .. |
| 2 | 1 | 1 |
| 2 | 1 | 2 |
| .. |  | .. |
| 6 | 6 | 6 |

- Instead of sending individual rolls, you roll the dice **three times**
  - And send the *triple* to your friend
- How many bits do you send *per roll?*
- 216 combinations: 8 bits per triple
  - Still 2.666 bits per roll
  - *Now we're talking!*

# Batching up 6-sided dice rolls

- Batching *four rolls*
  - 1296 combinations
  - 11 bits per outcome (4 rolls)
  - 2.75 bit per roll

- Batching *five rolls*
  - 7776 combinations
  - 13 bits per outcome (5 rolls)
  - 2.6 bits per roll

# Batching up 6-sided dice rolls



- Where will it end?

# Batching up 6-sided dice rolls



- Where will it end?

- $\lim_{k \to \infty} \frac{\lceil k \log2(6) \rceil}{k} = \log2(6)$ bits per roll in the limit

  – This is the absolute minimum – no batching will give you less than these many bits per outcome

# Can we do better?

- A four-sided die needs 2 bits per roll

- But then you find not all sides are equally likely

- P(1) = 0.5, P(2) = 0.25, P(3) 0.125, P(4) = 0.125

- *Can you do better than 2 bits per outcome*

# Can we do better?

- You have

P(1) = 0.5, P(2) = 0.25, P(3) 0.125, P(4) = 0.125

- You use:

| 1 | 0 |
|---|---|
| 2 | 1 0 |
| 3 | 1 1 0 |
| 4 | 1 1 1 |

  – Note receiver is *never in any doubt as to what they received*

- What is the average number of bits per outcome

# Can we do better?

- You have

P(1) = 0.5, P(2) = 0.25, P(3) 0.125, P(4) = 0.125

- You use:

| | |
|---|---|
| 1 | 0 |
| 2 | 1 0 |
| 3 | 1 1 0 |
| 4 | 1 1 1 |

  – Note receiver is *never in any doubt as to what they received*

- An outcome with probability $p$ is equivalent to obtaining one of $1/p$ equally likely choices

  – Requires $log2(\frac{1}{p})$ bits on average

# Entropy



- The average number of bits per symbol required to communicate a random variable over a digitial channel *using an optimal code* is

$$H(p) = \sum_i p_i \log \frac{1}{p_i} = -\sum_i p_i \log p_i$$

- You can't do better
  - Any other code will require more bits
- This is the *entropy of the random variable*

# 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

**X**

T,  M,  S…

M F  F M..

**Y**

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

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

# A brief review of basic info. theory



X → T, M, S…

Y → M F  F M..

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

- 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

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

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

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

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

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

# Onward..

# Projection: multiple notes

**M =**



**W =**



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

# We're actually computing a score

**M =**



H = ?

**W =**



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

# How about the other way?

M =



H =



W = ?

U = ?

- M ~ WH          W = Mpinv(H)        U = WH

# 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

$$\mathbf{W}, \mathbf{H} = \arg\min_{\overline{\mathbf{W}}, \overline{\mathbf{H}}} \| \mathbf{M} - \overline{\mathbf{W}\mathbf{H}} \|_F^2 + \Lambda(\overline{\mathbf{W}^T}\overline{\mathbf{W}} - \mathbf{I})$$

- Constraint: $\mathbf{W}$ is orthogonal
  - $\mathbf{W}^T\mathbf{W} = \mathbf{I}$
- The solution: $\mathbf{W}$ are the Eigen vectors of $\mathbf{MM^T}$
  - PCA!!

- $\mathbf{M} \sim \mathbf{WH}$ is an approximation
- Also, the rows of $\mathbf{H}$ are *decorrelated*
  - Trivial to prove that $\mathbf{HH^T}$ is diagonal

# PCA

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

$$\mathbf{M} \approx \mathbf{WH}$$

- The columns of W are the bases we have learned
  - The linear "building blocks" that compose the music
- They represent "learned" notes

# 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

# PCA through decorrelation of notes

$$\mathbf{W}, \mathbf{H} = \arg\min_{\overline{\mathbf{W}}, \overline{\mathbf{H}}} \| \mathbf{M} - \overline{\mathbf{H}} \|_F^2 + \Lambda(\overline{\mathbf{H}\mathbf{H}}^T - \mathbf{D})$$



- Different constraint: Constraint $\mathbf{H}$ to be decorrelated
  - $\mathbf{H}\mathbf{H}^T = \mathbf{D}$
- This will result exactly in PCA too
- Decorrelation of $\mathbf{H}$ Interpretation: What does this mean?

# 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..

# 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

- **Attempting to find statistically independent components of the mixed signal**
  - *Independent Component Analysis*

# Formulating it with Independence

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

- Impose statistical independence constraints on decomposition

# Changing problems for a bit



$$h_1(t)$$

$$m_1(t) = w_{11}h_1(t) + w_{12}h_2(t)$$

$$m_2(t) = w_{21}h_1(t) + w_{22}h_2(t)$$

$$h_2(t)$$

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

# A Separation Problem



- **M = WH**
  - **M** = "mixed" signal
  - **W** = "notes"
  - **H** = "transcription"

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

# A Separation Problem

**W**

$$\begin{matrix} w_{11} & w_{12} \\ w_{21} & w_{22} \end{matrix}$$

**H**



**M**



- Separation challenge: Given only $\mathbf{M}$ estimate $\mathbf{H}$

- **Identical to the problem of "finding notes"**

# Imposing Statistical Constraints

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



- **M = WH**

- Given only **M** estimate **H**

- $\mathbf{H} = \mathbf{W^{-1}M} = \mathbf{AM}$

- Only known constraint:  The rows of **H** are independent

- Estimate **A** such that the components of **AM** are statistically independent

  – **A** is the *unmixing* matrix

# Statistical Independence

- $M = WH$    $H = AM$ ← **Remember this form**

# An ugly algebraic solution

$$\mathbf{M} = \mathbf{WH} \quad \text{........} \quad \mathbf{H} = \mathbf{A}\mathbf{M}$$

- We could *decorrelate* signals by algebraic manipulation
  - We know uncorrelated signals have diagonal correlation matrix
  - So we transformed the signal so that it has a diagonal correlation matrix ($\mathbf{HH^T}$)

- Can we do the same for independence
  - Is there a linear transform that will enforce independence?

# An ugly algebraic solution

- We *decorrelated* signals by diagonalizing the covariance matrix

- *Is there a simple matrix we could just similarly diagonalize to make them independent?*

# An ugly algebraic solution

- We *decorrelated* signals by diagonalizing the covariance matrix

- *Is there a simple matrix we could just similarly diagonalize to make them independent?*
  - Not really, but there is a matrix we can diagonalize to make *fourth-order* moments independent
    - Just as decorrelation made second-order moments independent

# Emulating Independence

**H**



- 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

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

- If mean($\mathbf{M}$) = 0 => mean($\mathbf{H}$) = 0
  - $E[\mathbf{H}] = \mathbf{A}.E[\mathbf{M}] = \mathbf{A0} = \mathbf{0}$
  - First step of ICA: Set the mean of $\mathbf{M}$ to 0

$$\mu_{\mathbf{m}} = \frac{1}{cols\ (\mathbf{M})} \sum_{i} \mathbf{m}_i$$

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

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

# Emulating Independence..

H
H'
Diagonal
+ rank1 matrix

H=AM

A=BC

H=BCM

- Independence → Uncorrelatedness
- Find **C** such that **CM** is decorrelated
  - PCA
- Find **B** such that **B(CM)** is independent
- A little more than PCA

# Decorrelating and Whitening



$H$      $H'$     =    Diagonal   + rank1 matrix

$H = AM$

$A = BC$

$H = BCM$

- Eigen decomposition $\mathbf{MM}^T = \mathbf{ESE}^T$
- $\mathbf{C} = \mathbf{S}^{-1/2}\mathbf{E}^T$
- $\mathbf{X} = \mathbf{CM}$

- Not merely decorrelated but *whitened*
  - $\mathbf{XX}^T = \mathbf{CMM}^T\mathbf{C}^T = \mathbf{S}^{-1/2}\mathbf{E}^T\ \mathbf{ESE}^T\mathbf{ES}^{-1/2} = \mathbf{I}$

- $\mathbf{C}$ is the *whitening matrix*

# Uncorrelated != Independent

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

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

- 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

H

H'  $=$  Diagonal  $+$ **rank1**
**matrix**

H=AM

A=BC

H=BCM

- **X = CM**
- **XX$^T$ = I**

**H=BX**

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

# FOBI: Freeing Fourth Moments

- Find $\mathbf{B}$ such that the rows of $\mathbf{H} = \mathbf{BX}$ are independent

- The fourth moments of $\mathbf{H}$ have the form:
  $$\mathrm{E}[\mathbf{h}_i \, \mathbf{h}_j \, \mathbf{h}_k \, \mathbf{h}_l]$$

- If the rows of $\mathbf{H}$ were independent
  $$\mathrm{E}[\mathbf{h}_i \, \mathbf{h}_j \, \mathbf{h}_k \, \mathbf{h}_l] = \mathrm{E}[\mathbf{h}_i] \, \mathrm{E}[\mathbf{h}_j] \, \mathrm{E}[\mathbf{h}_k] \, \mathrm{E}[\mathbf{h}_l]$$

- Solution: Compute $\mathbf{B}$ such that the fourth moments of $\mathbf{H} = \mathbf{BX}$ are decoupled
  - While ensuring that $\mathbf{B}$ is Unitary

- **FOBI: Fourth Order Blind Identification**

# ICA: Freeing Fourth Moments

$$\mathbf{H} =$$



$h_k$

Objective: Find a matrix B such that the rows of H=BX are statistically independent

Define a matrix D that *would* be diagonal if the rows of BX are independent

Compute B such that this matrix becomes diagonal

- Create a matrix of fourth moment terms that would be diagonal were the rows of $\mathbf{H}$ independent and diagonalize it

- A good candidate: the weighted correlation matrix of $\mathbf{H}$

$$\boldsymbol{D} = E\left[\|\boldsymbol{h}\|^2 \boldsymbol{h}\boldsymbol{h}^{\mathrm{T}}\right] = \sum_k \|\boldsymbol{h}_k\|^2 \boldsymbol{h}_k \boldsymbol{h}_k^{\mathrm{T}}$$

  – $\boldsymbol{h}$ are the columns of $\mathbf{H}$
  – Assuming $\boldsymbol{h}$ is real, else replace transposition with Hermitian

# ICA: The D matrix

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

$$\boldsymbol{D} = \boldsymbol{E}\left[\|\boldsymbol{h}\|^2 \boldsymbol{h}\boldsymbol{h}^{\mathrm{T}}\right]$$

$$d_{ij} = \boldsymbol{E}\left[\left(\sum_l h_l^2\right) h_i h_j\right]$$

Sum of squares of all components

$$\sum_l h_l^2$$

$i^{th}$ component

$j^{th}$ component

$$h_i \, h_j$$

$$\left(\sum_l h_l^2\right) h_i h_j$$

## On the actual matrix

$$\boldsymbol{D} = \sum_k \|\boldsymbol{h}_k\|^2 \boldsymbol{h}_k \boldsymbol{h}_k^{\mathrm{T}}$$

$$d_{ij} = \frac{1}{cols(\mathbf{H})} \sum_k \left(\sum_l h_{kl}^2\right) h_{ki} h_{kj}$$

# ICA: The D matrix

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

$$\boxed{D = E[\|h\|^2 hh^{\mathrm{T}}]} \qquad d_{ij} = E\left[\left(\sum_l h_l^2\right) h_i h_j\right]$$

$$d_{ij} = \frac{1}{cols(\mathbf{H})} \sum_k \left(\sum_l h_{kl}^2\right) h_{ki} h_{kj}$$

- If the $h_i$ terms were independent and zero mean
- For $i\ != j$

$$E\left[h_i h_j \sum_l h_l^2\right] = E[h_i^3]E[h_j] + E[h_i]E[h_j^3] + E[h_i]E[h_j] \sum_{l \neq i, l \neq j} E[h_l^3] = 0$$

- For $i = j$

  - $E\left[h_i h_j \sum_l h_l^2\right] = E[h_i^4] + E[h_i^2] \sum_{l \neq i} E[h_l^2] \neq 0$

- i.e., if $h_i$ were independent, $D$ would be a diagonal matrix
  - **Let us diagonalize $D$**

84

# Diagonalizing D

- Recall: $\mathbf{H} = \mathbf{B}\mathbf{X}$
  - $\mathbf{B}$ is what we're trying to learn to make $\mathbf{H}$ independent
  - Assumption: $\mathbf{B}$ is unitary, i.e. $\mathbf{B}\mathbf{B}^T = \mathbf{I}$

- Note: if $\mathbf{H} = \mathbf{B}\mathbf{X}$, then each vector $\mathbf{h} = \mathbf{B}\mathbf{x}$
- The fourth moment matrix of $\mathbf{H}$ is
- $\mathbf{D} = \mathrm{E}[\mathbf{h}^T \, \mathbf{h} \, \mathbf{h} \, \mathbf{h}^T] = \mathrm{E}[\mathbf{x}^T\mathbf{B}\mathbf{B}^T\mathbf{x} \, \mathbf{B}^T \, \mathbf{x} \, \mathbf{x}^T\mathbf{B}]$

$$= \mathrm{E}[\mathbf{x}^T\mathbf{x} \, \mathbf{B}^T \, \mathbf{x} \, \mathbf{x}^T\mathbf{B}]$$

$$= \mathbf{B}^T \, \mathrm{E}[\mathbf{x}^T\mathbf{x} \, \mathbf{x}\mathbf{x}^T]\mathbf{B}$$

$$= \mathbf{B}^T \, \mathrm{E}[\|\mathbf{x}\|^2 \, \mathbf{x}\mathbf{x}^T]\mathbf{B}$$

Objective: Find a matrix B such that the rows of H=BX are statistically independent

Define a matrix D that *would* be diagonal if the rows of BX are independent

Compute B such that this matrix becomes diagonal

# Diagonalizing D

- Objective: Estimate $\mathbf{B}$ such that the fourth moment of $\mathbf{H} = \mathbf{BX}$ is diagonal

- Compose $\mathbf{D_x} = \sum_k \|\mathbf{x}_k\|^2 \mathbf{x_k} \mathbf{x}_k^{\mathrm{T}}$

- Diagonalize $\mathbf{D}_x$ via Eigen decomposition

  $\mathbf{D}_x = \mathbf{U} \Lambda \mathbf{U}^{\mathrm{T}}$

- $\mathbf{B} = \mathbf{U}^{\mathrm{T}}$
  - **That's it!!!!**

# B frees the fourth moment

$$\mathbf{D_x} = \mathbf{U}\Lambda\mathbf{U}^T \ ; \ \mathbf{B} = \mathbf{U}^T$$

- $\mathbf{U}$ is a unitary matrix, i.e. $\mathbf{U}^T\mathbf{U} = \mathbf{U}\mathbf{U}^T = \mathbf{I}$ (identity)

- $\mathbf{H} = \mathbf{BX} = \mathbf{U}^T\mathbf{X}$

- $\mathbf{h} = \mathbf{U}^T\mathbf{x}$

- The fourth moment matrix of $\mathbf{H}$ is

$$\mathrm{E}[\|\mathbf{h}\|^2 \ \mathbf{h} \ \mathbf{h}^T] = \mathbf{U}^T \ \mathrm{E}[\|\mathbf{x}\|^2 \ \mathbf{xx}^T]\mathbf{U}$$
$$= \mathbf{U}^T \ \mathbf{D_x} \ \mathbf{U}$$
$$= \mathbf{U}^T \ \mathbf{U} \ \Lambda \ \mathbf{U}^T \ \mathbf{U} = \Lambda$$

- The fourth moment matrix of $\mathbf{H} = \mathbf{U}^T\mathbf{X}$ is Diagonal!!

# Overall Solution

- Objective: Estimate A such that the rows of $\mathbf{H} =$ $\mathbf{AM}$ are independent

- Step 1: *Whiten M*
  - $\mathbf{C}$ is the (transpose of the) matrix of Eigen vectors of $\mathbf{MM}^{\mathrm{T}}$
  - $\mathbf{X} = \mathbf{CM}$

- Step 2: Free up fourth moments on $\mathbf{X}$
  - $\mathbf{B}$ is the (transpose of the) matrix of Eigenvectors of $\mathbf{X}.diag(\mathbf{X}^{\mathrm{T}}\mathbf{X}).\mathbf{X}^{\mathrm{T}}$
  - $\mathbf{A} = \mathbf{BC}$

# FOBI for ICA

- Goal: to derive a matrix **A** such that the rows of **AM** are independent

- Procedure:

  1. "Center" **M**

  2. Compute the autocorrelation matrix $\boldsymbol{R}_{MM}$ of **M**

  3. Compute whitening matrix **C** via Eigen decomposition
     $$\boldsymbol{R}_{MM} = \mathbf{ESE}^T, \quad \mathbf{C} = \mathbf{S}^{-1/2}\mathbf{E}^T$$

  4. Compute $\mathbf{X} = \mathbf{CM}$

  5. Compute the fourth moment matrix $\mathbf{D}' = E[\|\mathbf{x}\|^2 \mathbf{xx}^T]$

  6. Diagonalize **D'** via Eigen decomposition

  7. $\mathbf{D}' = \mathbf{U\Lambda U}^T$

  8. Compute $\mathbf{A} = \mathbf{U}^T \boldsymbol{C}$

- The fourth moment matrix of **H**=**AM** is diagonal

  – **Note that the autocorrelation matrix of H will also be diagonal**

# ICA by diagonalizing moment matrices

- FOBI is not perfect
  - Only a subset of fourth order moments are considered
    - 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 multiple fourth-order cumulant matrices

# Enforcing Independence

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

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

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

- Define and minimize a contrast function

  » F($\mathbf{AM}$)

- Contrast functions are often only *approximations* too..

# A note on pre-whitening

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

- Eigen decomposition $\mathbf{MM}^T = \mathbf{ESE}^T$
- $\mathbf{C} = \mathbf{S}^{-1/2}\mathbf{E}^T$

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

- $\mathbf{X} = \mathbf{CM}$
  - $E[\mathbf{x}_i\mathbf{x}_j] = \delta_{ij}$ 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 : $\mathbf{H} = \mathbf{BX}$
  - $\mathbf{X}$ is "whitened" $\mathbf{M}$

# Linear Functions

- $\mathbf{h} = \mathbf{B}\mathbf{x}, \quad \mathbf{x} = \mathbf{B}^{-1}\mathbf{h}$

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

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

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

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

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

# 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

- Recall PCA

- $\mathbf{M} = \mathbf{WH}$, the columns of $\mathbf{W}$ must be orthogonal

- Leads to: $\min_{\mathbf{W}} ||\mathbf{M} - \mathbf{WW}^{\mathrm{T}}\mathbf{M}||^2 + \Lambda.\mathrm{trace}(\mathbf{W}^{\mathrm{T}}\mathbf{W})$
  - Error minimization framework to estimate $\mathbf{W}$

- Can we arrive at an error minimization framework for ICA

- Define an "Error" objective that represents independence

# An alternate approach

- Definition of Independence – if $x$ and $y$ are independent:
  - $\mathrm{E}[f(x)g(y)] = \mathrm{E}[f(x)]\mathrm{E}[g(y)]$
  - 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

- $\mathbf{P} = \mathbf{g(H)} \, \mathbf{f(H)}^{\mathrm{T}} = \mathbf{g(BX)} \, \mathbf{f(BX)}^{\mathrm{T}}$

$$\mathbf{P} = \begin{bmatrix} P_{11} & P_{21} & \dots \\ P_{12} & P_{22} & \\ \cdot & \cdot & \\ \cdot & \cdot & \\ \cdot & \cdot & \end{bmatrix}$$

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

This is a square matrix

- Must ideally be

$$\mathbf{Q} = \begin{bmatrix} Q_{11} & Q_{21} & \dots \\ Q_{12} & Q_{22} & \\ \cdot & \cdot & \\ \cdot & \cdot & \\ \cdot & \cdot & \end{bmatrix}$$

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

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

- Error $= \|\mathbf{P}\text{-}\mathbf{Q}\|_{\mathrm{F}}^{2}$

# An alternate approach

- Ideal value for $\mathbf{Q}$

$$\mathbf{Q} = \begin{bmatrix} Q_{11} & Q_{21} & \cdots \\ Q_{12} & Q_{22} \\ \cdot & \cdot \\ \cdot & \cdot \\ \cdot & \cdot \end{bmatrix}$$

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

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

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

# An alternate approach

- Minimize Error

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

$$\mathbf{Q} = Diagonal$$

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

- Leads to trivial Widrow Hopf type iterative rule:

$$\mathbf{E} = Diag - \mathbf{g(BX)f(BX)^T}$$

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

# Update Rules

- Multiple solutions under different assumptions for g() and f()

- $\mathbf{H} = \mathbf{BX}$

- $\mathbf{B} = \mathbf{B} + \eta\,\Delta\mathbf{B}$

- Jutten Herraut : Online update

  - $\Delta B_{ij} = f(\mathbf{h}_i)g(\mathbf{h}_j)$;  -- actually assumed a recursive neural network

- Bell Sejnowski

  - $\Delta\mathbf{B} = ([\mathbf{B}^T]^{-1} - \mathbf{g(H)X}^T)$

# Update Rules

- Multiple solutions under different assumptions for g() and f()

- $\mathbf{H} = \mathbf{BX}$

- $\mathbf{B} = \mathbf{B} + \eta\,\Delta\mathbf{B}$

- Natural gradient  -- f() = identity function
  - $\Delta\mathbf{B} = (\mathbf{I} - \mathbf{g}(\mathbf{H})\mathbf{H}^{\mathrm{T}})\,\mathbf{X}^{\mathrm{T}}$

- Cichoki-Unbehaeven
  - $\Delta\mathbf{B} = (\mathbf{I} - \mathbf{g}(\mathbf{H})\mathbf{f}(\mathbf{H})^{\mathrm{T}})\,\mathbf{X}^{\mathrm{T}}$

# What are G() and F()

- 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
  - $\Delta \mathbf{B} = (\mathbf{I} - \mathbf{HH}^T\text{-}\mathbf{Ktanh(H)H}^T)\ \mathbf{X}^T$
- Or simply
  - $\Delta \mathbf{B} = (\mathbf{I} - \mathbf{Ktanh(H)H}^T)\ \mathbf{X}^T$

# 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*



PCA
ICA

# 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 Enhncement



- Very commonly used to enhance EEG signals
- EEG signals are frequently corrupted by heartbeats and biorhythm signals
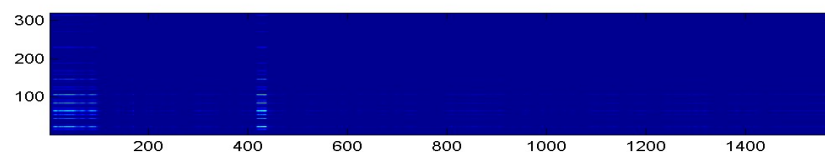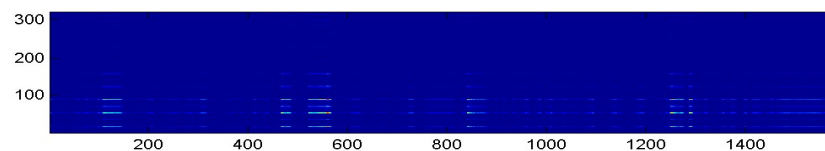- ICA can be used to separate them out
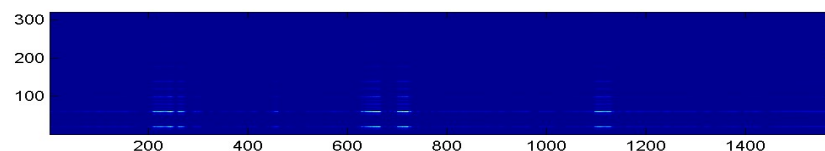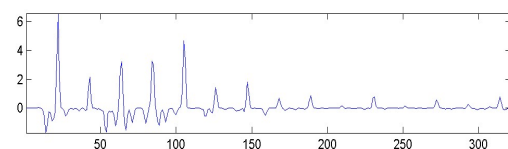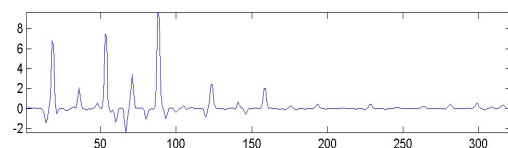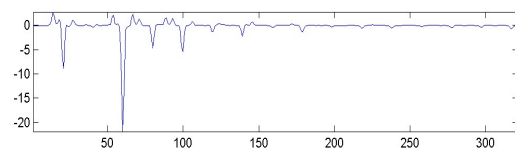
# So how does that work?



- 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..