

ENGR-E 511; ENGR-E 399

“Machine Learning for Signal Processing”

Module 04: Dimension Reduction

Minje Kim

Department of Intelligent Systems Engineering

Email: minje@indiana.edu

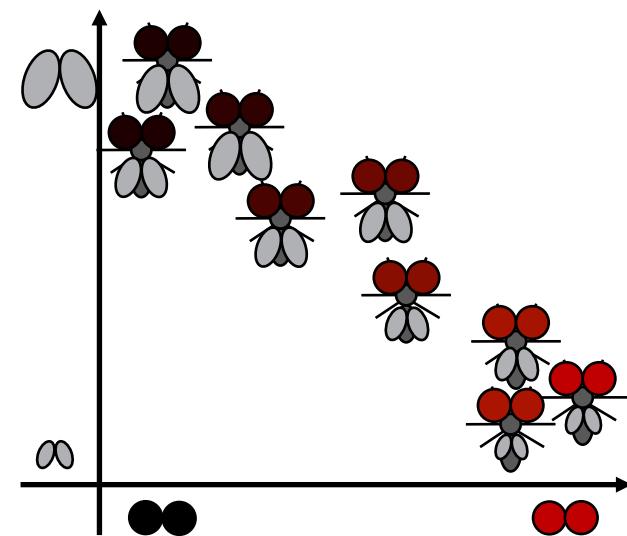
Website: <http://minjekim.com>

Research Group: <http://saige.sice.indiana.edu>

Meeting Request: <http://doodle.com/minje>

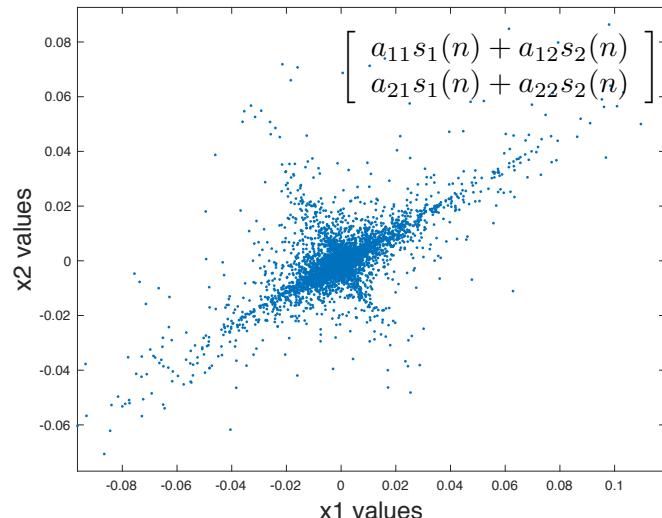
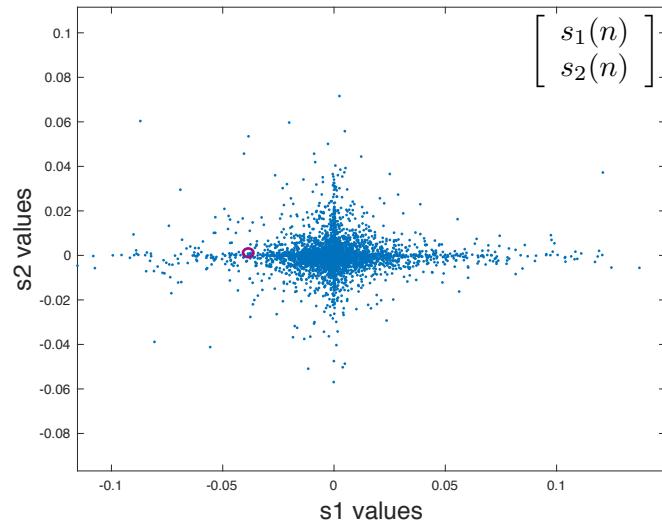


INDIANA UNIVERSITY
**SCHOOL OF INFORMATICS,
COMPUTING, AND ENGINEERING**



A Toy Source Separation Problem

- Two sensors and two sources
 - Two people are saying at the same time and two mics are recording it
 - I'm ignoring a lot of important things
 - e.g. delay, reverberation, geometry, etc.
 - The sources
 - s_1 : 
 - s_2 : 
 - We can scatter plot these samples on a 2D space (a $2 \times N$ matrix)
 - The projection onto each axis recovers sources
 - Each recording is a linear mixture of the sources
 - x_1 : 
 - x_2 : 
 - Now, the projection onto each axis doesn't recover the sources
 - We need to project the data points more carefully
 - What would those correct projections be?

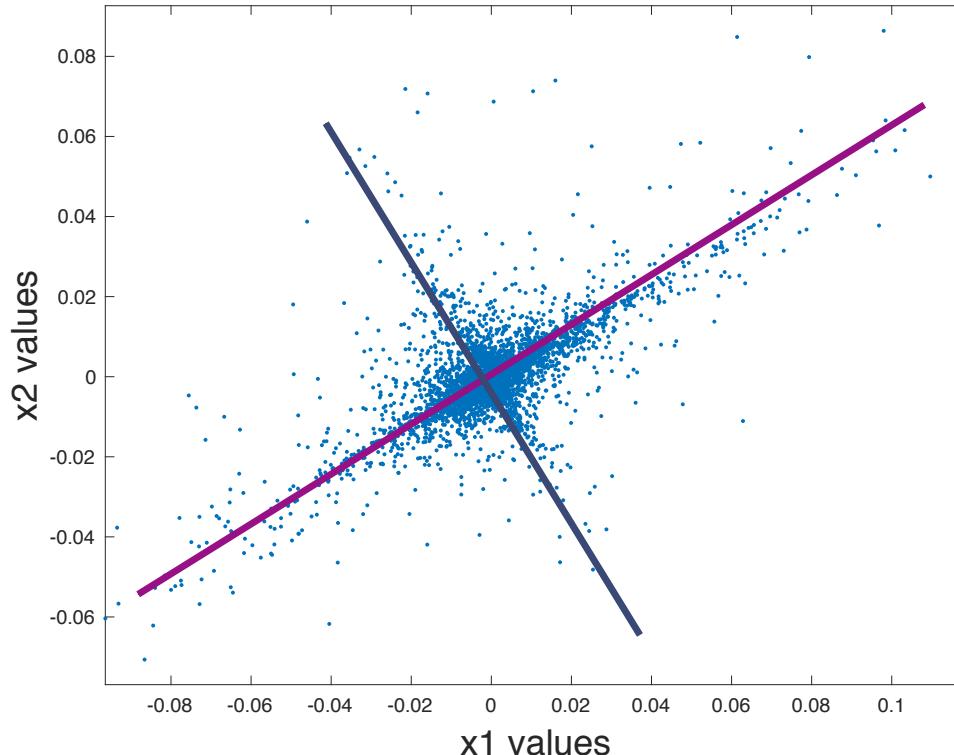


INDIANA UNIVERSITY

SCHOOL OF INFORMATICS, COMPUTING, AND ENGINEERING

A Toy Source Separation Problem

- Two sensors and two sources
 - For this overly simplified example, we know the answer
 - Because the linear combination for mixing is nothing but some rotation



- How do we find these vectors?



INDIANA UNIVERSITY

SCHOOL OF INFORMATICS, COMPUTING, AND ENGINEERING

Principal Component Analysis

- a.k.a. Karhunen-Loève Transform

- First, let's define the variables $\mathbf{x}_n = \begin{bmatrix} x_1(n) \\ x_2(n) \end{bmatrix} = \begin{bmatrix} A_{11}s_1(n) + A_{12}s_2(n) \\ A_{21}s_1(n) + A_{22}s_2(n) \end{bmatrix} = \mathbf{A}\mathbf{s}_n$

$$\mathbf{X} = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N]$$

- For this data I've got a feeling that the best projection would be

 - The one that makes the sample variance maximal

- If we assume zero mean variables, the variance after a projection is

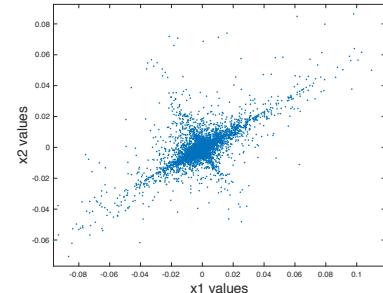
$$var(\mathbf{w}^\top \mathbf{x}_n) = \frac{1}{N} \sum_n (\mathbf{w}^\top \mathbf{x}_n)^2 = \frac{1}{N} \sum_n \left([w_1, w_2] \begin{bmatrix} x_1(n) \\ x_2(n) \end{bmatrix} \right)^2$$

 - Or in the matrix form $= \frac{1}{N} (\mathbf{w}^\top \mathbf{X})(\mathbf{w}^\top \mathbf{X})^\top = \frac{1}{N} \mathbf{w}^\top \mathbf{X} \mathbf{X}^\top \mathbf{w}$

- The objective function is $\arg \max_{\mathbf{w}} \frac{1}{N} \mathbf{w}^\top \mathbf{X} \mathbf{X}^\top \mathbf{w}$

 - What am I missing? $\arg \max_{\mathbf{w}} \mathbf{w}^\top \Sigma_x \mathbf{w} + \lambda(1 - \mathbf{w}^\top \mathbf{w})$

 - What happens if it were not for the constraints?



INDIANA UNIVERSITY

SCHOOL OF INFORMATICS, COMPUTING, AND ENGINEERING

Principal Component Analysis

- a.k.a. Karhunen-Loève Transform

- Partial differentiation $\arg \max_{\mathbf{w}} \mathbf{w}^\top \Sigma_x \mathbf{w} + \lambda(1 - \mathbf{w}^\top \mathbf{w})$

$$\frac{\partial \mathcal{J}}{\partial \mathbf{w}} = 2\Sigma_x \mathbf{w} - 2\lambda \mathbf{w} = 0 \quad \frac{\partial \mathcal{J}}{\partial \lambda} = 1 - \mathbf{w}^\top \mathbf{w} = 0 \quad \Sigma_x \mathbf{w} = \lambda \mathbf{w}$$

- It's an eigendecomposition problem
- The variance after projection equals the eigenvalue

$$\mathbf{w}^\top \Sigma_x \mathbf{w} = \lambda \Leftrightarrow \mathbf{w}^\top \mathbf{X} \mathbf{X}^\top \mathbf{w} = \lambda$$

$$\mathbf{W}^\top \Sigma_x \mathbf{W} = \Lambda$$

- Principal Components (PC)
 - They are projection vectors
 - Can be found using eigendecomposition on the covariance matrix (eigenvectors)
 - They are ordered based on the eigenvalues
 - \mathbf{w}_1 is the 1st PC if $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_D$



INDIANA UNIVERSITY

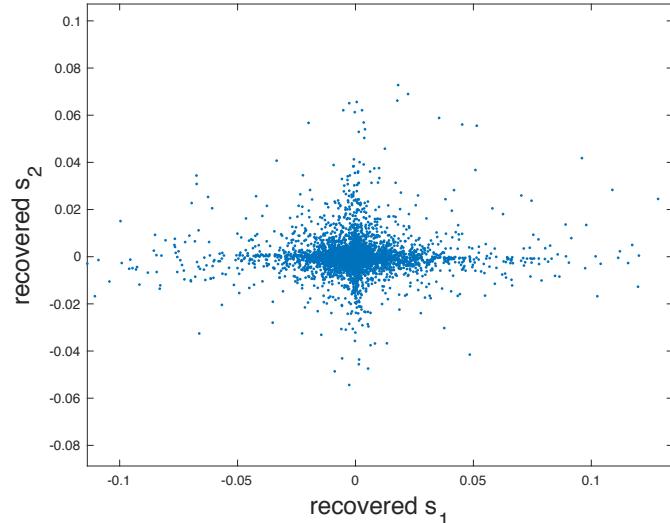
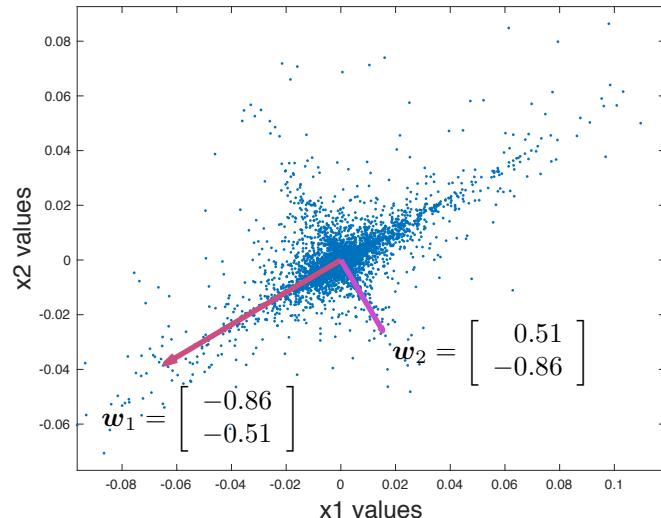
SCHOOL OF INFORMATICS, COMPUTING, AND ENGINEERING

Principal Component Analysis

- a.k.a. Karhunen-Loève Transform

- The estimated PCs
 - There's a sign ambiguity
- Variances after projection
 - $\lambda_1 = 3.48 \times 10^{-4}$, $\lambda_2 = 5.65 \times 10^{-5}$
- How about source separation?
 - We're trying to eliminate the mixing effect by rotating the data points back to the original space
 - Therefore: $\mathbf{W}^\top \mathbf{A} \approx \mathbf{I}$ $\therefore \mathbf{x}_n = \mathbf{A}\mathbf{s}_n$
 - $$\begin{bmatrix} -0.86 & -0.51 \\ 0.51 & -0.86 \end{bmatrix} \begin{bmatrix} 0.87 & -0.50 \\ 0.50 & 0.86 \end{bmatrix} = \begin{bmatrix} -1 & -0.01 \\ 0.01 & -1 \end{bmatrix}$$
 - Almost correct except the flipped signs
- Separation (projection) results

$$\begin{bmatrix} s_1 \\ s_2 \end{bmatrix} \approx \mathbf{W}^\top \mathbf{X}$$



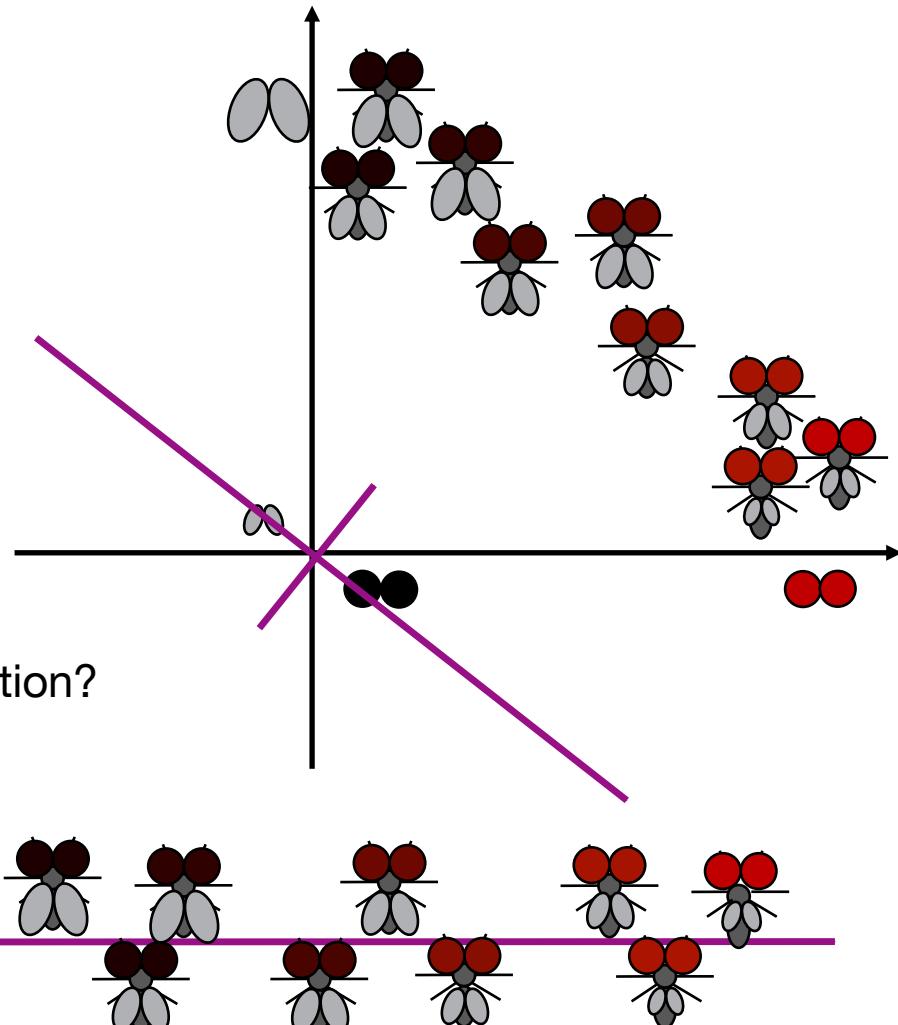
INDIANA UNIVERSITY

SCHOOL OF INFORMATICS, COMPUTING, AND ENGINEERING

PCA for Dimension Reduction

- The Mating Flies and Lazy RA Example

- One day as he takes off Prof. K ordered his RA to group the flies into two clusters
 - The RA forgot about it and went home..
- Originally it was a clustering problem
 - GMM could have described the sample distribution
 - But, not anymore, b/c over the night the flies mated
- Is there another way to describe this sample distribution?
 - Rather than clustering?
 - The 1st PC
- We lost some information by reducing the dimension?



INDIANA UNIVERSITY

SCHOOL OF INFORMATICS, COMPUTING, AND ENGINEERING

Clustering vs PCA

- In terms of compression

- Clustering (GMM, k-means, vector quantization, etc)
 - It reduces the levels of representation by allowing “collision”
 - Multiple data samples are allowed (or forced) to have the same value
 - Number of clusters will define the amount of compression
- PCA (or the other dimension reduction methods)
 - It transforms the data into a new coordinate system
 - Which can be defined with a less number of axis
 - Each data sample has its own unique combination of coefficients in the new coordinate system (no collision)
 - Still need to encode the coefficients



INDIANA UNIVERSITY

SCHOOL OF INFORMATICS, COMPUTING, AND ENGINEERING

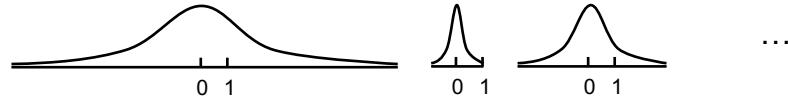
PCA for Whitening

- Feature extraction and normalization in one shot

- Say that there are 10 original dimensions with unique dimension-wise distributions

- Our job is

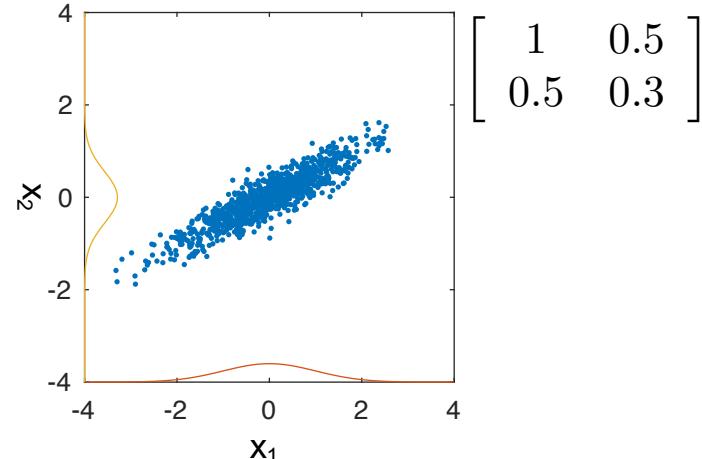
- To figure out a smaller set of features
 - To lessen the difference in dynamic ranges



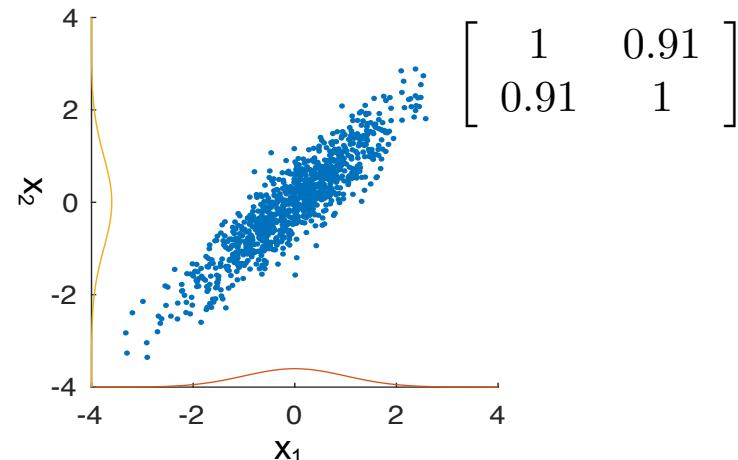
- As for the DR, maybe just a standardization could be enough

- Zero mean and unit variance

- But, some of them are correlated



After coordinate-wise standardization



INDIANA UNIVERSITY

SCHOOL OF INFORMATICS, COMPUTING, AND ENGINEERING

PCA for Whitening

- Feature extraction and normalization in one shot
 - Whitening is not only for ensuring unit variances (1's in the diagonal)
 - But for making off-diagonals zero
 - We can use PCA for whitening
 - Eigendecomposition of the covariance yields a diagonal matrix: $\mathbf{W}^\top \Sigma_x \mathbf{W} = \Lambda$ ← So far, just diagonalization, not whitening
 - By the way,
- $$\Lambda = \Lambda^{\frac{1}{2}} \Lambda^{\frac{1}{2}} = \begin{bmatrix} \lambda_1^{\frac{1}{2}} & 0 & \cdots & 0 \\ 0 & \lambda_2^{\frac{1}{2}} & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \lambda_D^{\frac{1}{2}} \end{bmatrix} \begin{bmatrix} \lambda_1^{\frac{1}{2}} & 0 & \cdots & 0 \\ 0 & \lambda_2^{\frac{1}{2}} & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \lambda_D^{\frac{1}{2}} \end{bmatrix}$$
- Therefore,
- $$\Lambda^{(-\frac{1}{2})} \mathbf{W}^\top \Sigma_x \mathbf{W} = \Lambda^{\frac{1}{2}} \quad \Lambda^{(-\frac{1}{2})} \mathbf{W}^\top \Sigma_x \mathbf{W} \Lambda^{(-\frac{1}{2})} = \mathbf{I}$$
- The projection for whitening:
- $$\Lambda^{(-\frac{1}{2})} \mathbf{W}^\top \mathbf{X} = \begin{bmatrix} \mathbf{w}_1^\top / \sqrt{\lambda_1} \\ \mathbf{w}_2^\top / \sqrt{\lambda_2} \\ \vdots \\ \mathbf{w}_D^\top / \sqrt{\lambda_D} \end{bmatrix} \begin{bmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \\ \vdots \\ \mathbf{x}_D \end{bmatrix}$$

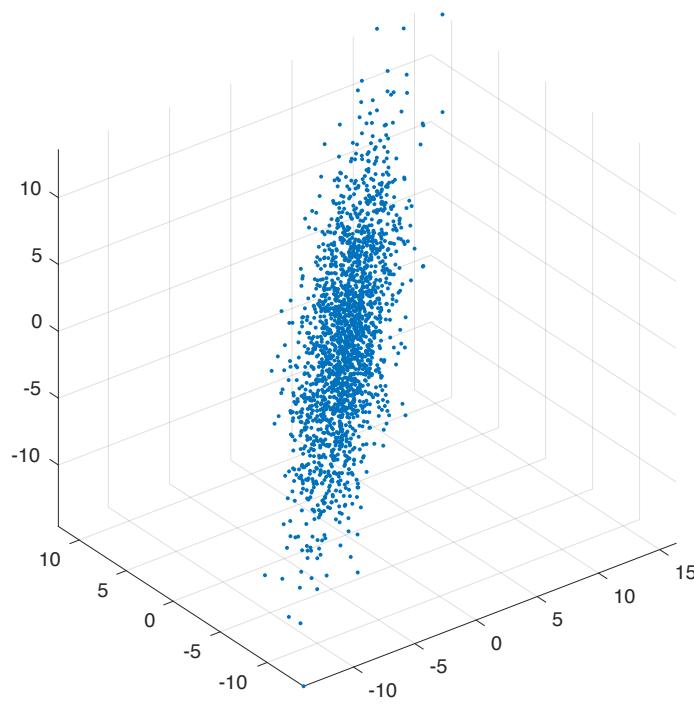


INDIANA UNIVERSITY

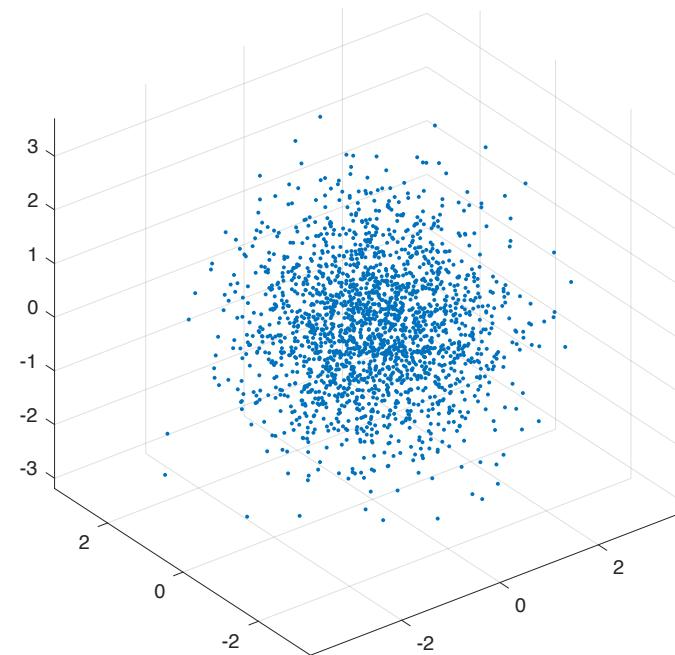
SCHOOL OF INFORMATICS, COMPUTING, AND ENGINEERING

PCA for Whitening

- Feature extraction and normalization in one shot
 - A whitening example (MATLAB fig)



Before whitening



After whitening



INDIANA UNIVERSITY

SCHOOL OF INFORMATICS, COMPUTING, AND ENGINEERING

Electroencephalography (EEG) Signals

- A Higher-Dimensional Data

- Our first audio example in this lecture wasn't really a dimension reduction case
 - It's more likely a rotating problem for unmixing
- An EEG data with 14 channels and a label for the eye status {open, close}

```
@RELATION EEG_DATA
@ATTRIBUTE AF3 NUMERIC
@ATTRIBUTE F7 NUMERIC
@ATTRIBUTE F3 NUMERIC
@ATTRIBUTE FC5 NUMERIC
@ATTRIBUTE T7 NUMERIC
@ATTRIBUTE P7 NUMERIC
@ATTRIBUTE O1 NUMERIC
@ATTRIBUTE O2 NUMERIC
@ATTRIBUTE P8 NUMERIC
@ATTRIBUTE T8 NUMERIC
@ATTRIBUTE FC6 NUMERIC
@ATTRIBUTE F4 NUMERIC
@ATTRIBUTE F8 NUMERIC
@ATTRIBUTE AF4 NUMERIC
@ATTRIBUTE eyeDetection {0,1}
@DATA
```

```
4329.23,4009.23,4289.23,4148.21,4350.26,4586.15,4096.92,4641.03,4222.05,4238.46,4211.28,4280.51,4635.90,4393.85,0
4324.62,4004.62,4293.85,4148.72,4342.05,4586.67,4097.44,4638.97,4210.77,4226.67,4207.69,4279.49,4632.82,4384.10,0
4327.69,4006.67,4295.38,4156.41,4336.92,4583.59,4096.92,4630.26,4207.69,4222.05,4206.67,4282.05,4628.72,4389.23,0
4328.72,4011.79,4296.41,4155.90,4343.59,4582.56,4097.44,4630.77,4217.44,4235.38,4210.77,4287.69,4632.31,4396.41,0
4362.56,3983.08,4273.33,4115.38,4336.41,4624.62,4104.62,4627.18,4212.82,4233.85,4245.13,4324.62,4663.08,4421.03,1
4358.46,3978.97,4268.21,4104.62,4330.26,4615.90,4103.08,4623.08,4205.13,4235.90,4244.62,4322.05,4663.08,4411.79,1
4363.08,3981.03,4276.92,4106.67,4333.85,4608.21,4097.44,4614.36,4194.36,4230.26,4247.69,4323.59,4667.18,4418.46,1
4366.67,3981.54,4274.36,4112.82,4334.87,4612.82,4093.85,4608.72,4192.82,4226.67,4247.18,4318.97,4671.28,4426.67,1
...
```



<https://en.wikipedia.org/wiki/Electroencephalography>



<https://en.wikipedia.org/wiki/Emotiv>



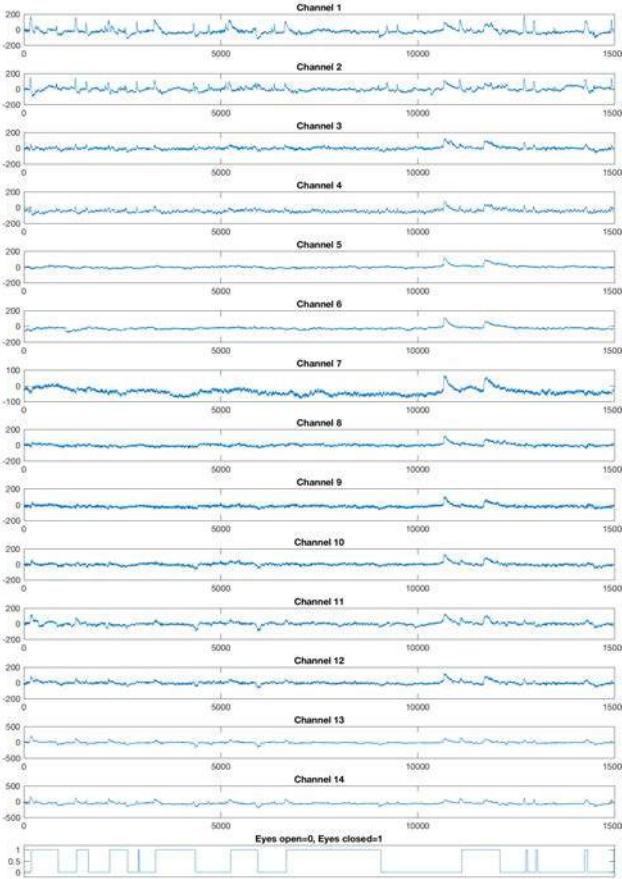
INDIANA UNIVERSITY

SCHOOL OF INFORMATICS, COMPUTING, AND ENGINEERING

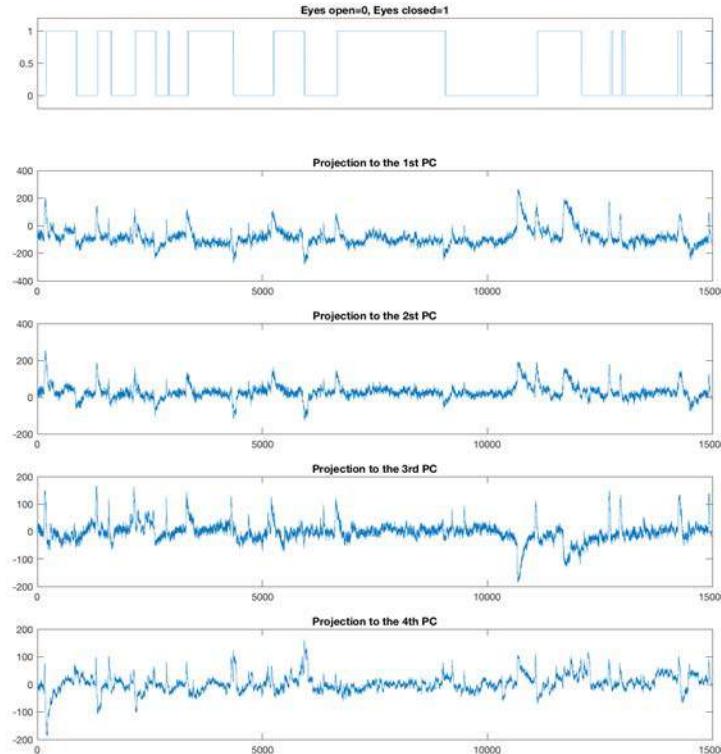
<https://archive.ics.uci.edu/ml/datasets/EEG+Eye+State>

Electroencephalography (EEG) Signals

- A Higher-Dimensional Data



Raw signals (after outlier detection)



After PCA (4 PCs)



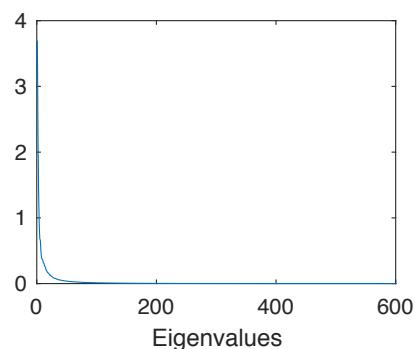
INDIANA UNIVERSITY

SCHOOL OF INFORMATICS, COMPUTING, AND ENGINEERING

Eigenfaces

- Just another dimension reduction example

- Dimension reduction procedure



INDIANA UNIVERSITY

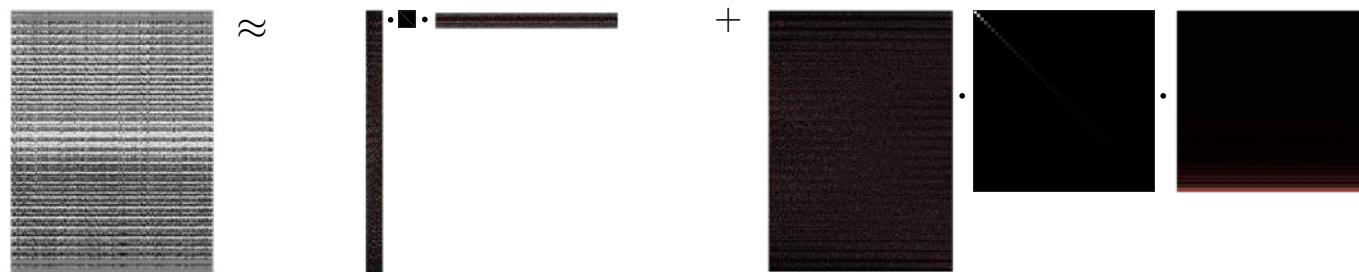
SCHOOL OF INFORMATICS, COMPUTING, AND ENGINEERING

Eigenfaces

- Just another dimension reduction example

- Why does it work?
- Remember SVD
 - whose singular vectors are eigenvectors of the covariance matrix

$$\mathbf{X} \approx \mathbf{V} \mathbf{S} \mathbf{U}^\top = \underbrace{\mathbf{V}_{:,1:50} \mathbf{S}_{1:50,1:50} \mathbf{U}_{:,1:50}^\top}_{\text{More significant part}} + \underbrace{\mathbf{V}_{:,51:D} \mathbf{S}_{51:D,51:D} \mathbf{U}_{:,51:D}^\top}_{\text{Less significant part}}$$



- So, the dimension reduction projections are

$$\mathbf{W}_{:,1:50}^\top \mathbf{X} = \mathbf{V}_{:,1:50}^\top \mathbf{X} \approx \mathbf{V}_{:,1:50}^\top \mathbf{V}_{:,1:50} \mathbf{S}_{1:50,1:50} \mathbf{U}_{:,1:50}^\top + 0 \xrightarrow{\text{Orthogonality}}$$

$$\mathbf{W}_{:,1:50}^\top \mathbf{X} = \mathbf{S}_{1:50,1:50} \mathbf{U}_{:,1:50}^\top = \mathbf{Z}_{1:50,:} \xleftarrow{\text{Unnormalized right singular vectors}}$$

- Then, the reconstruction procedure is

$$\mathbf{W}_{:,1:50} \mathbf{Z}_{1:50,:} = \mathbf{W}_{:,1:50} \mathbf{W}_{:,1:50}^\top \mathbf{X} = \mathbf{V}_{:,1:50} \mathbf{S}_{1:50,1:50} \mathbf{U}_{:,1:50}^\top \xleftarrow{\text{SVD with first 50 components}}$$

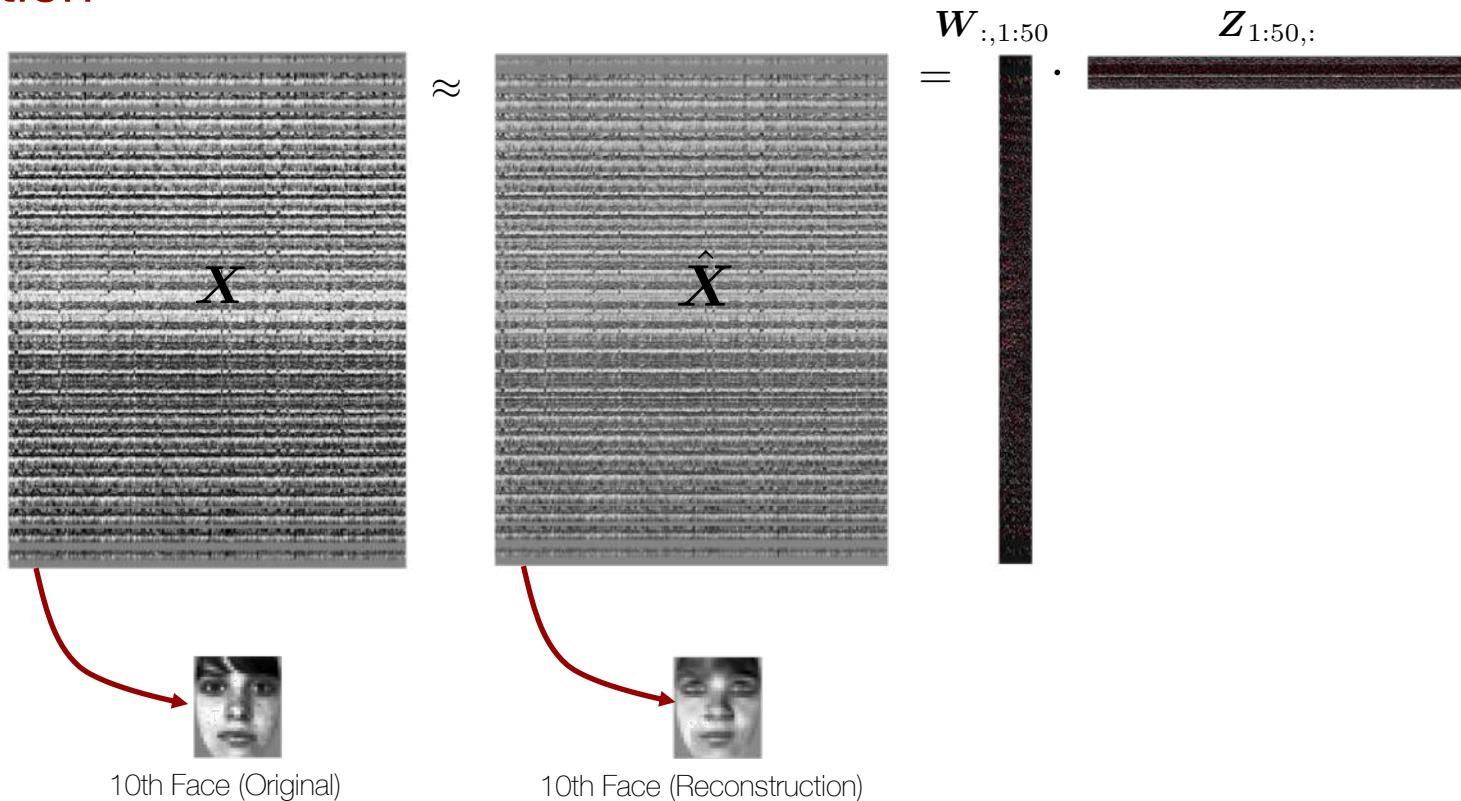


INDIANA UNIVERSITY

SCHOOL OF INFORMATICS, COMPUTING, AND ENGINEERING

Eigenfaces

- Reconstruction



- We see that only 50 coefficients can carry a lot of information



INDIANA UNIVERSITY

SCHOOL OF INFORMATICS, COMPUTING, AND ENGINEERING

Eigenfaces

- The learned eigenfaces (eigenvectors in the matrix form)



INDIANA UNIVERSITY

SCHOOL OF INFORMATICS, COMPUTING, AND ENGINEERING

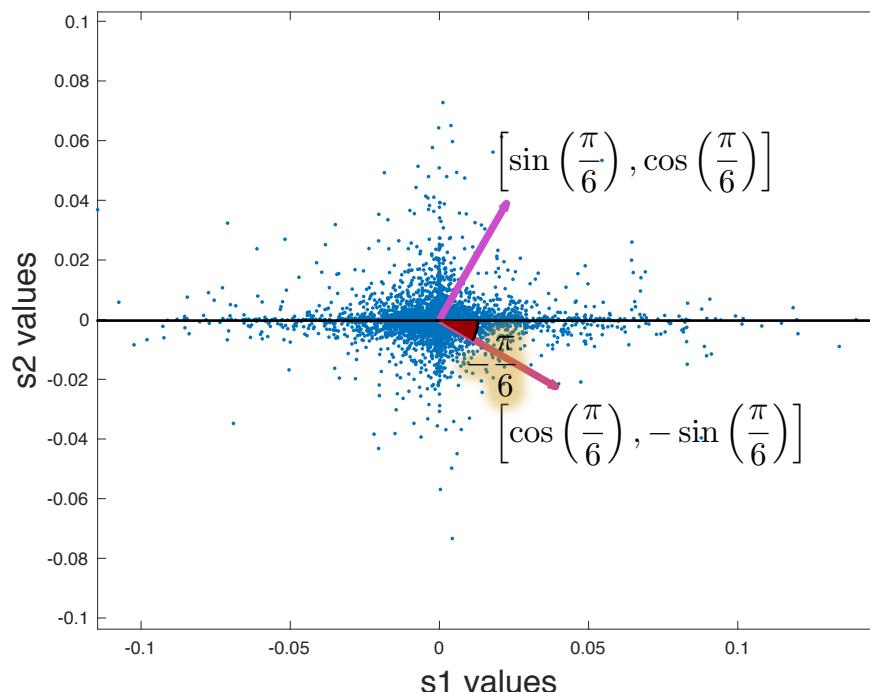
Another Toy Source Separation Problem

- A little bit more realistic example
 - In the previous separation example, the mixing matrix I used was
 - It's a special matrix, because the two rotations are orthogonal

$$\left[\cos\left(\frac{\pi}{6}\right), -\sin\left(\frac{\pi}{6}\right) \right] \cdot \left[\begin{array}{c} \sin\left(\frac{\pi}{6}\right) \\ \cos\left(\frac{\pi}{6}\right) \end{array} \right] = 0$$

- So, the separation was exceptionally easy
 - Just rotating the data back to normal
- What if the mixing matrix was more realistic?

$$A = \begin{bmatrix} \cos\left(\frac{\pi}{6}\right) & -\sin\left(\frac{\pi}{6}\right) \\ \sin\left(\frac{\pi}{6}\right) & \cos\left(\frac{\pi}{6}\right) \end{bmatrix}$$



INDIANA UNIVERSITY

SCHOOL OF INFORMATICS, COMPUTING, AND ENGINEERING

Another Toy Source Separation Problem

- A little bit more realistic example

- A more realistic mixing matrix

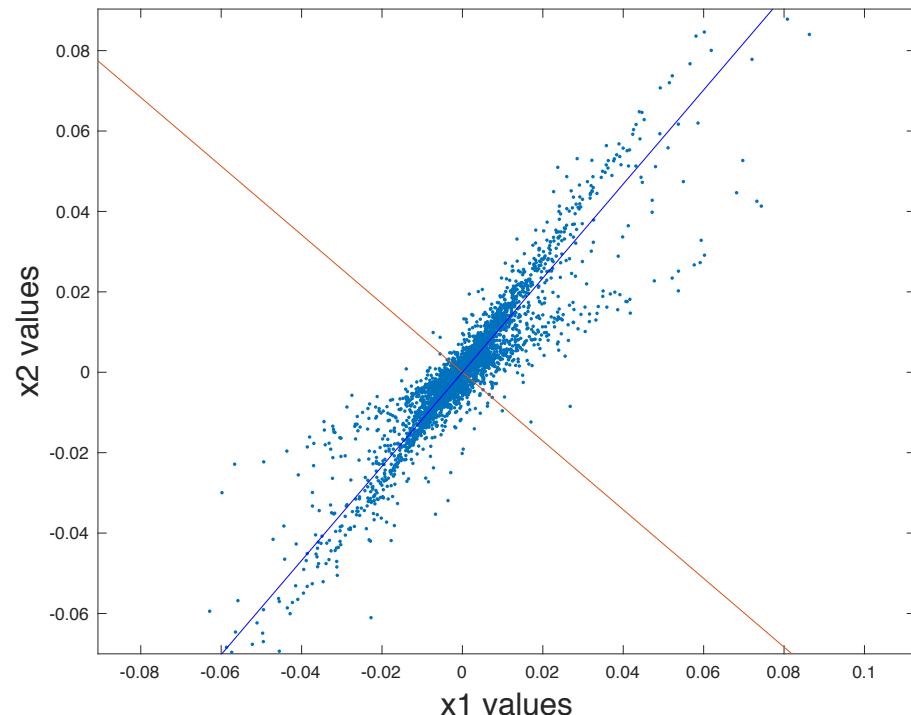
$$\mathbf{A} = \begin{bmatrix} 0.5 & 0.9 \\ 0.7 & 0.4 \end{bmatrix}, \quad \begin{bmatrix} x_1(n) \\ x_2(n) \end{bmatrix} = \begin{bmatrix} 0.5 & 0.9 \\ 0.7 & 0.4 \end{bmatrix} \cdot \begin{bmatrix} s_1(n) \\ s_2(n) \end{bmatrix}$$

- What does it mean?

- In the first channel the female speaker is attenuated more (0.5 vs 0.9)
 - In the second channel the male speaker is attenuated more (0.4 vs 0.7)

- PCA?

- Doesn't seem to work



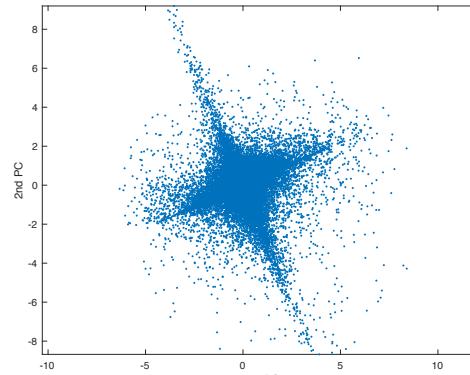
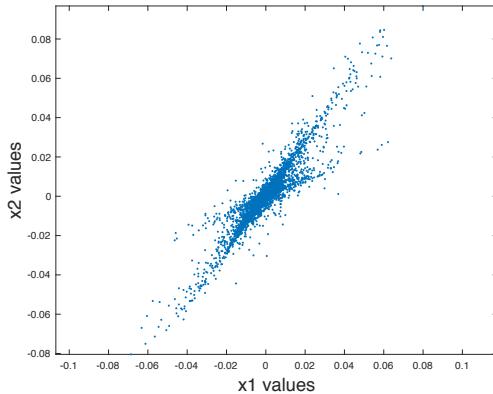
INDIANA UNIVERSITY

SCHOOL OF INFORMATICS, COMPUTING, AND ENGINEERING

Another Toy Source Separation Problem

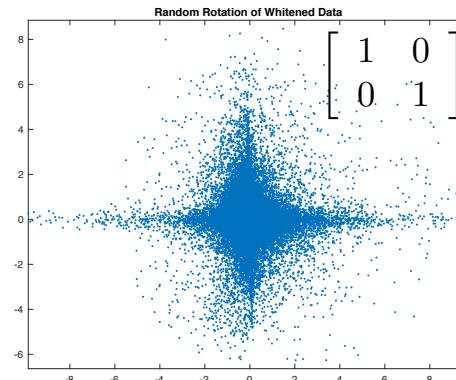
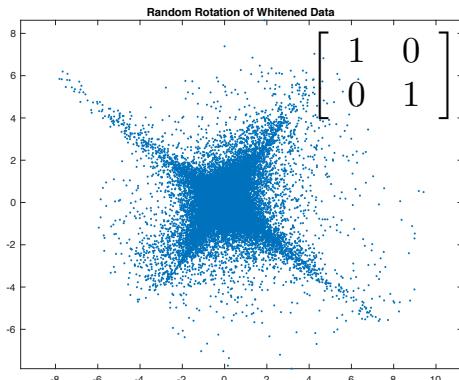
- A little bit more realistic example

- Can PCA whitening solve this?



$$\text{cov}(\mathbf{Z}) = \mathbf{I}$$

- No, because after whitening all rotations give unit covariance



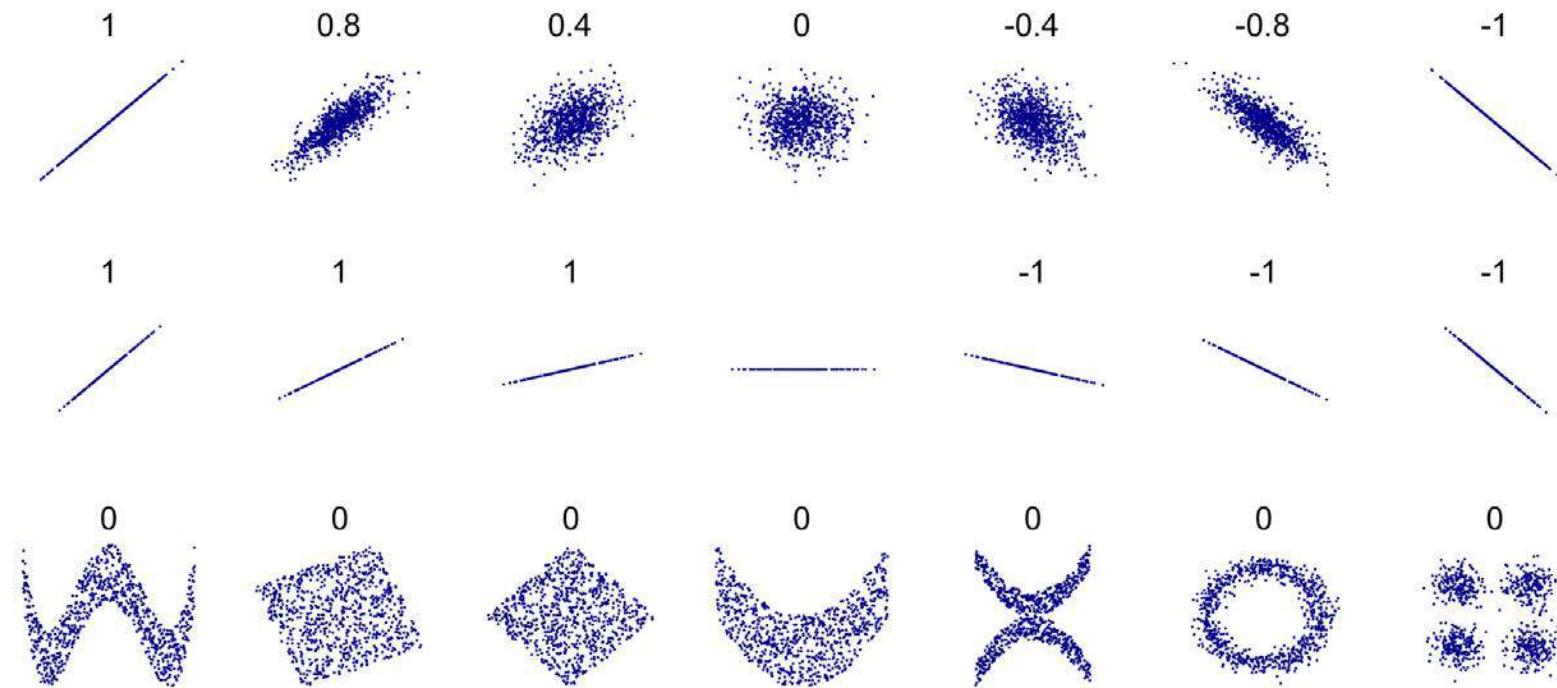
INDIANA UNIVERSITY

SCHOOL OF INFORMATICS, COMPUTING, AND ENGINEERING

Correlation

- Whitening is not enough
 - Correlation values of various joint distributions

$$\frac{cov(x, y)}{\sigma_x \sigma_y}$$



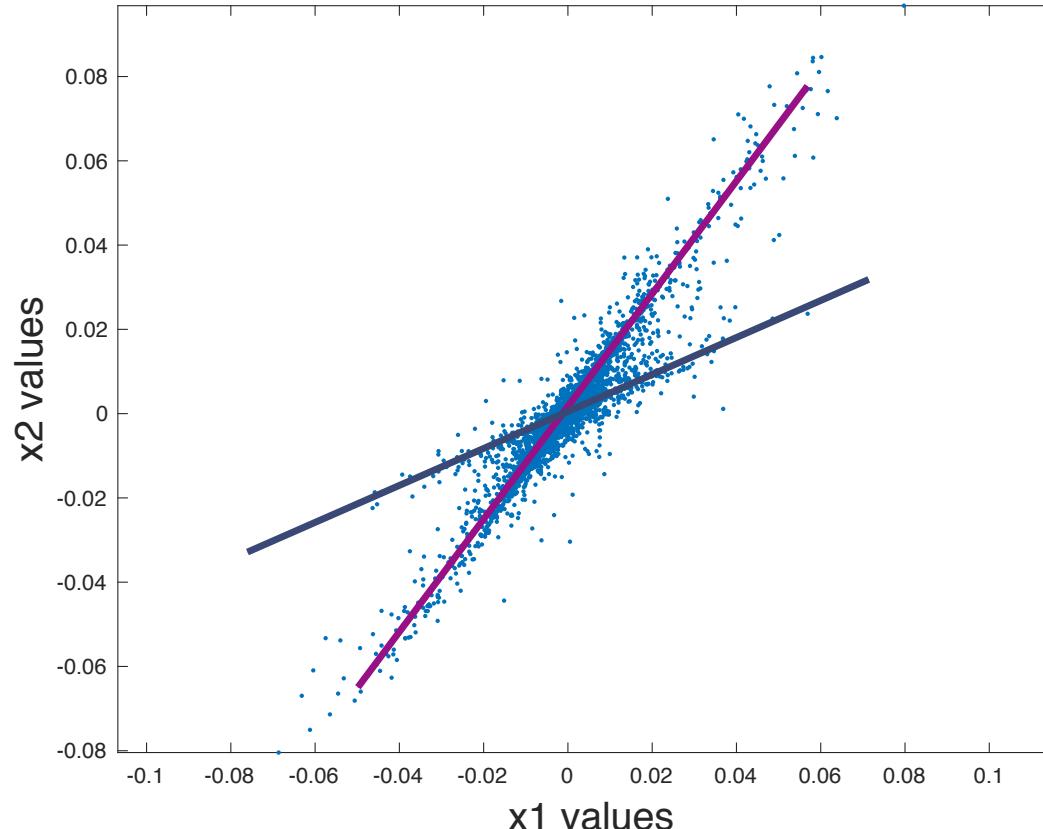
INDIANA UNIVERSITY

SCHOOL OF INFORMATICS, COMPUTING, AND ENGINEERING

https://en.wikipedia.org/wiki/Correlation_and_dependence

Independent Component Analysis

- An algorithm beyond whitening
- We need another criteria that can estimate the non-orthogonal basis vectors



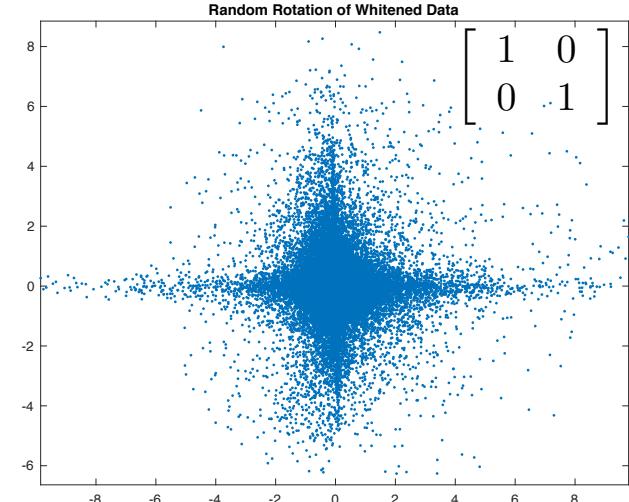
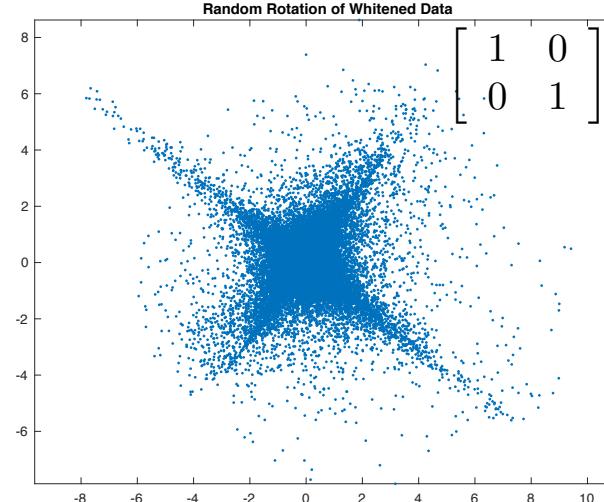
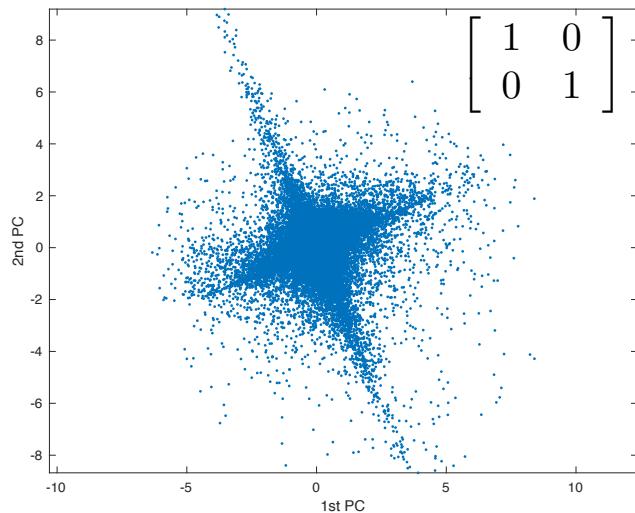
INDIANA UNIVERSITY

SCHOOL OF INFORMATICS, COMPUTING, AND ENGINEERING

Independent Component Analysis

- An algorithm beyond whitening

- All these whitened data and their rotations are already uncorrelated $E[x, y] = E[x]E[y]$



- How about their independence? $P(x, y) = P(x)P(y)$

$$E[f(x)g(y)] = \int_x \int_y f(x)g(y)P(x)P(y)dxdy = \int_x f(x)P(x)dx \int_y g(y)P(y)dy = E[f(x)]E[g(y)]$$

- In order to be independent, for any functions f and g , $f(x)$ and $g(y)$ should be uncorrelated
- This is a much stronger condition than whitening

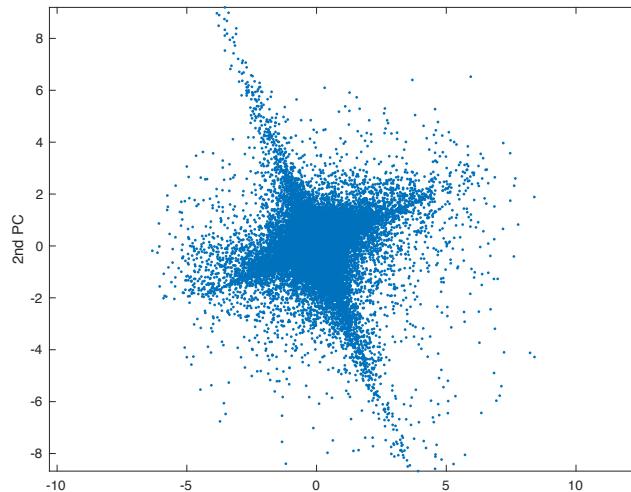


INDIANA UNIVERSITY

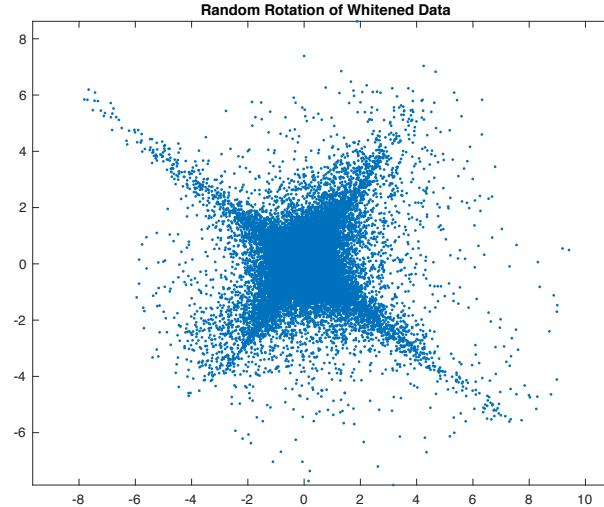
SCHOOL OF INFORMATICS, COMPUTING, AND ENGINEERING

Independent Component Analysis

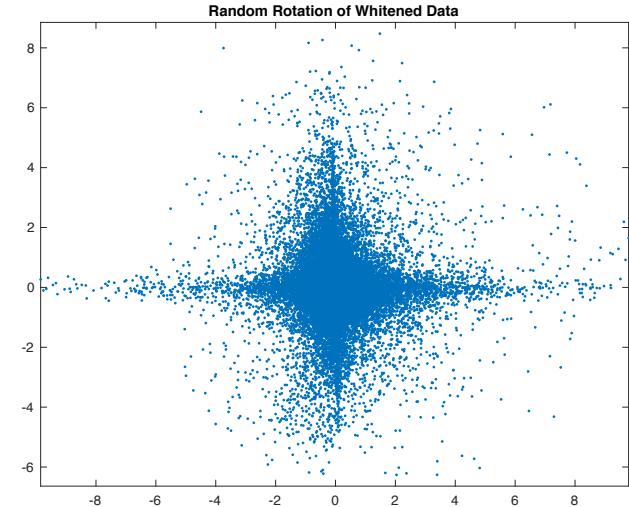
- An algorithm beyond whitening
- Let's check the covariance matrix of after applying $f(x)$ and $g(y)$



$$\begin{aligned} \text{cov}(\tanh(x), \tanh(y)) &= 0.0157 \\ \text{cov}(x^3, y^3) &= -42.63 \end{aligned}$$



$$\begin{aligned} 0.0154 \\ -69.68 \end{aligned}$$



$$\begin{aligned} -0.0059 \\ 8.55 \end{aligned}$$

- There's a rotation that can maximize the independence
- A Usual pipeline for ICA
 1. Whiten the data using PCA
 1. Drop some eigenvectors if dimension reduction is needed
 2. Do ICA



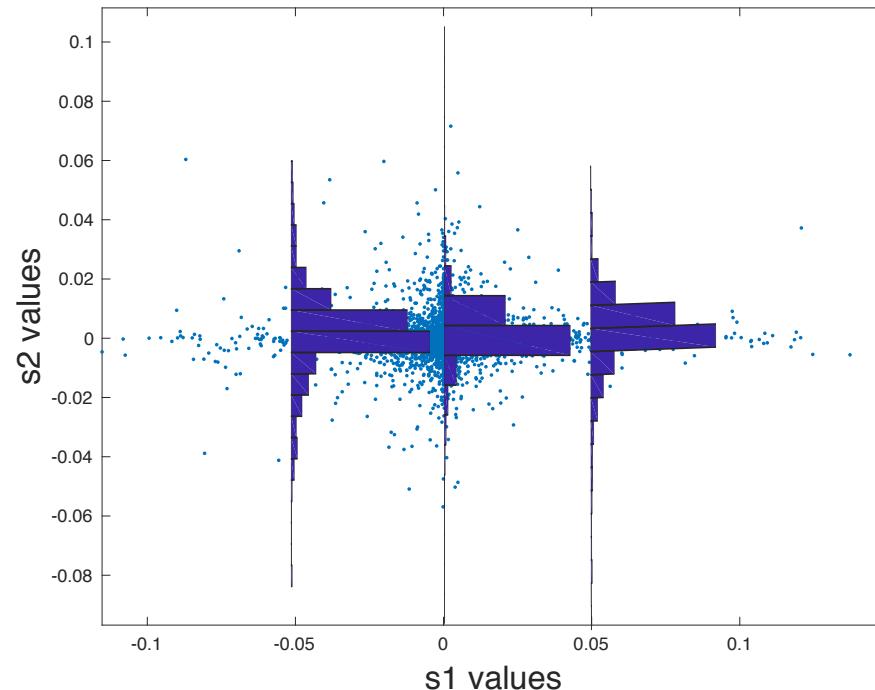
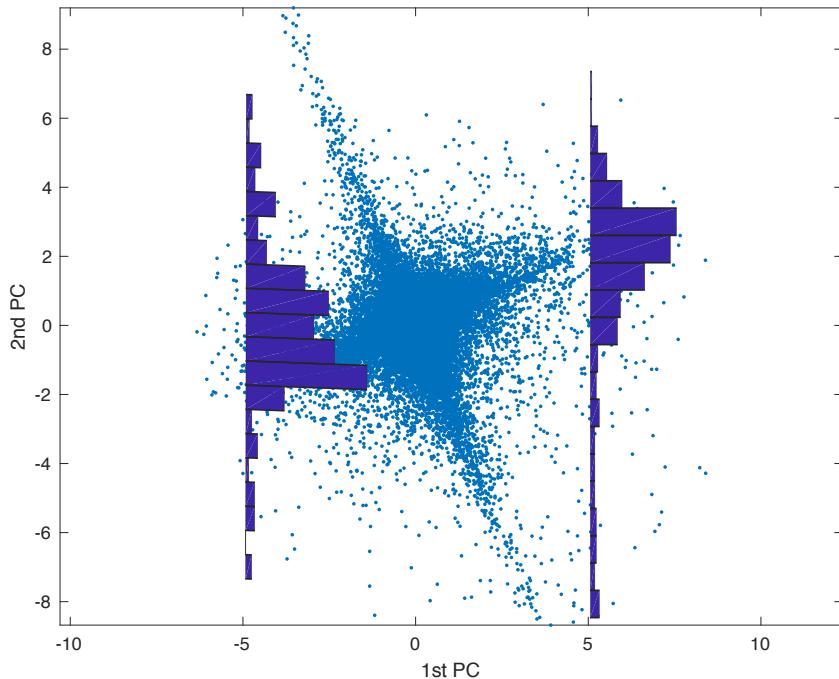
INDIANA UNIVERSITY

SCHOOL OF INFORMATICS, COMPUTING, AND ENGINEERING

Independent Component Analysis

- An algorithm beyond whitening

- Which one is more independent? $P(x_2|x_1) = P(x_2)$?



INDIANA UNIVERSITY

SCHOOL OF INFORMATICS, COMPUTING, AND ENGINEERING

Independent Component Analysis

- The optimization

- There are a lot of ICA algorithms. A successful ICA can produce
 - The components that are less Gaussian-like (Central Limit Theorem)
 - Higher-order statistics can measure the non-Gaussianity
 - The components whose **joint dist** and **product of marginal distributions** are with minimal KL-div
 - According to the definition of independence
- One update rule derived from the mutual information

$$\Delta \mathbf{W} \propto (DI - f(\mathbf{y}) \cdot g(\mathbf{y})^\top) \mathbf{W}$$

\mathbf{W} : ICA unmixing matrix

$$\mathbf{W} \leftarrow \mathbf{W} + \rho \Delta \mathbf{W}$$

\mathbf{y} : Estimated sources

D : The number of data samples

$f(\cdot), g(\cdot)$: Nonlinear functions

ρ : learning rate

- In what condition $\Delta \mathbf{W} = 0$?

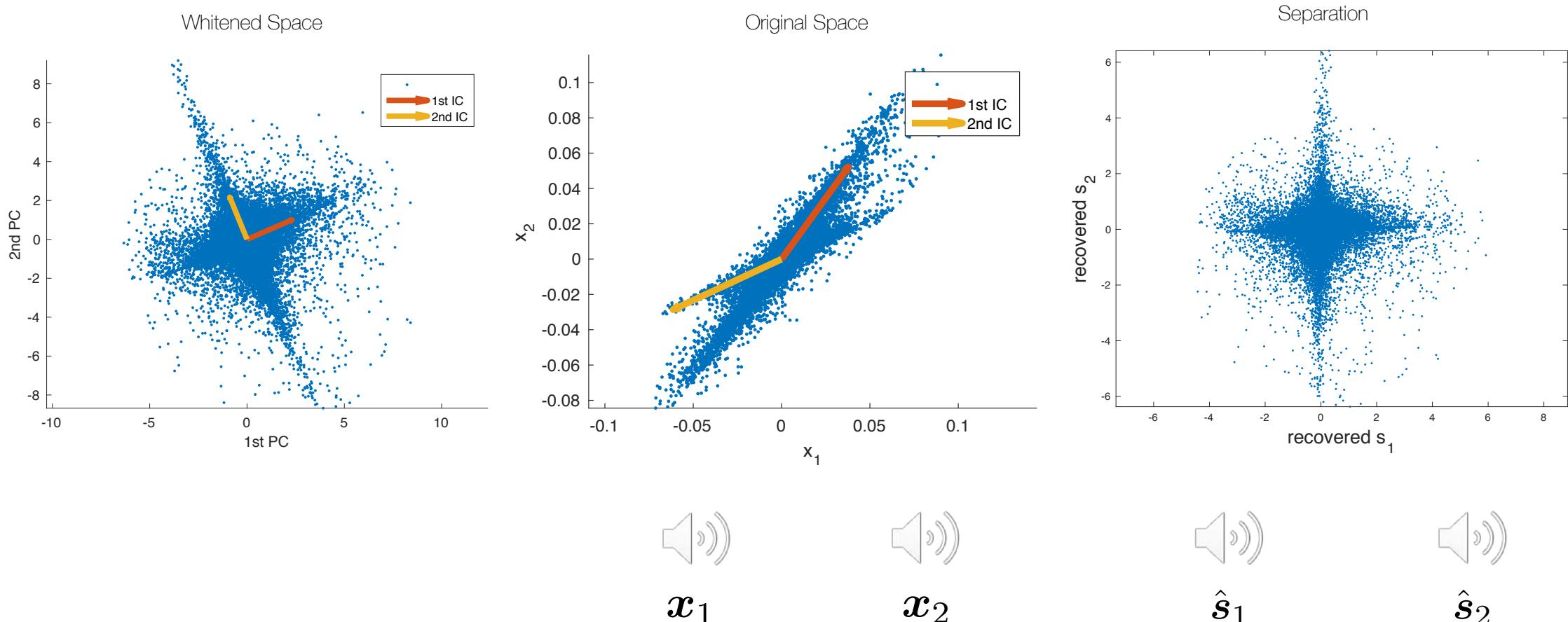


INDIANA UNIVERSITY

SCHOOL OF INFORMATICS, COMPUTING, AND ENGINEERING

ICA for Separation

- Distribution of the recovered sources



INDIANA UNIVERSITY

SCHOOL OF INFORMATICS, COMPUTING, AND ENGINEERING

Properties of ICA

- Identifiability Condition

- ICA wants to find an inverse of the mixing matrix \mathbf{A} to recover the sources (kind of indirectly)

$$\mathbf{X} = \mathbf{AS}$$

$$\mathbf{Z} = \underbrace{(\Lambda^{-1/2} \mathbf{W}_{pca}^\top)}_{\text{Whitening}} \mathbf{X} = \underbrace{(\Lambda^{-1/2} \mathbf{W}_{pca}^\top)}_{\text{Whitening}} \mathbf{AS}$$

$$\mathbf{W}_{ica} \mathbf{Z} = \mathbf{W}_{ica} (\Lambda^{-1/2} \mathbf{W}_{pca}^\top) \mathbf{X} = \underbrace{\mathbf{W}_{ica} (\Lambda^{-1/2} \mathbf{W}_{pca}^\top)}_{\text{Left Inverse of } \mathbf{A}} \mathbf{AS}$$

- \mathbf{A} needs to have more rows than columns (or at least square) and be full column rank
 - Sources need to be non-Gaussians

- ICA solutions are subject to ambiguities

- A permuted and scaled version of the sources:

$$\begin{bmatrix} 0 & .3 \\ 2 & 0 \end{bmatrix} \mathbf{S}$$

- We can build a new mixing scenario, from which ICA starts:

$$\mathbf{X} = \mathbf{A} \underbrace{\begin{bmatrix} 0 & .5 \\ 1/0.3 & 0 \end{bmatrix}}_{\text{The new mixing matrix}} \begin{bmatrix} 0 & .3 \\ 2 & 0 \end{bmatrix} \mathbf{S}$$

- There are infinitely many permuted/scaled versions of the mixing matrix



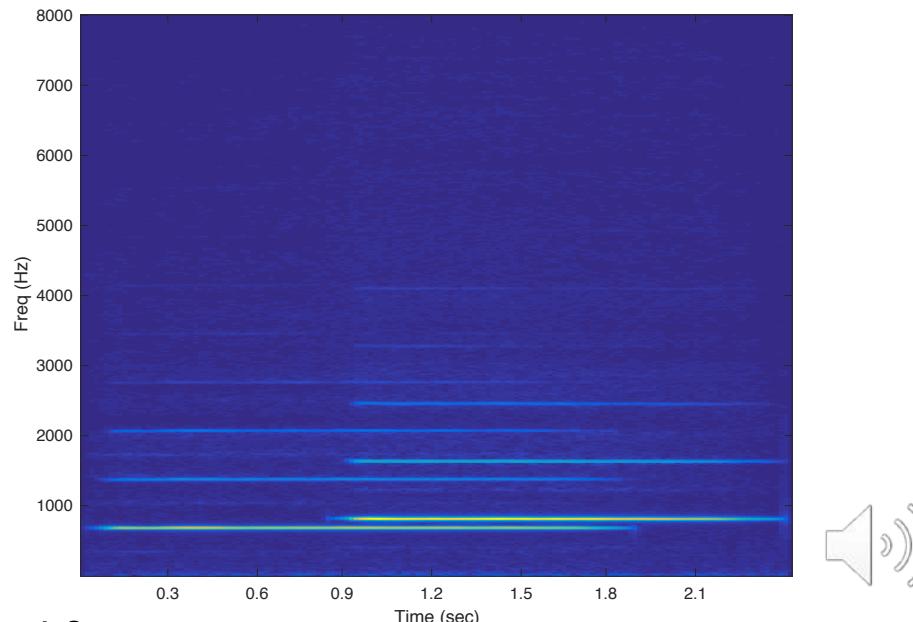
INDIANA UNIVERSITY

SCHOOL OF INFORMATICS, COMPUTING, AND ENGINEERING

Underdetermined Separation

- What if we have only one mixture
 - What if the mixing matrix A is just a row vector?
 - Or, what if the mixing procedure is not linear?
 - A single-channel recording with two sources

$$x(n) = [a_1, a_2] \begin{bmatrix} s_1(n) \\ s_2(n) \end{bmatrix}$$



- Would a clustering technique work?
 - No (overlap of sources from 0.9 to 1.8 sec)

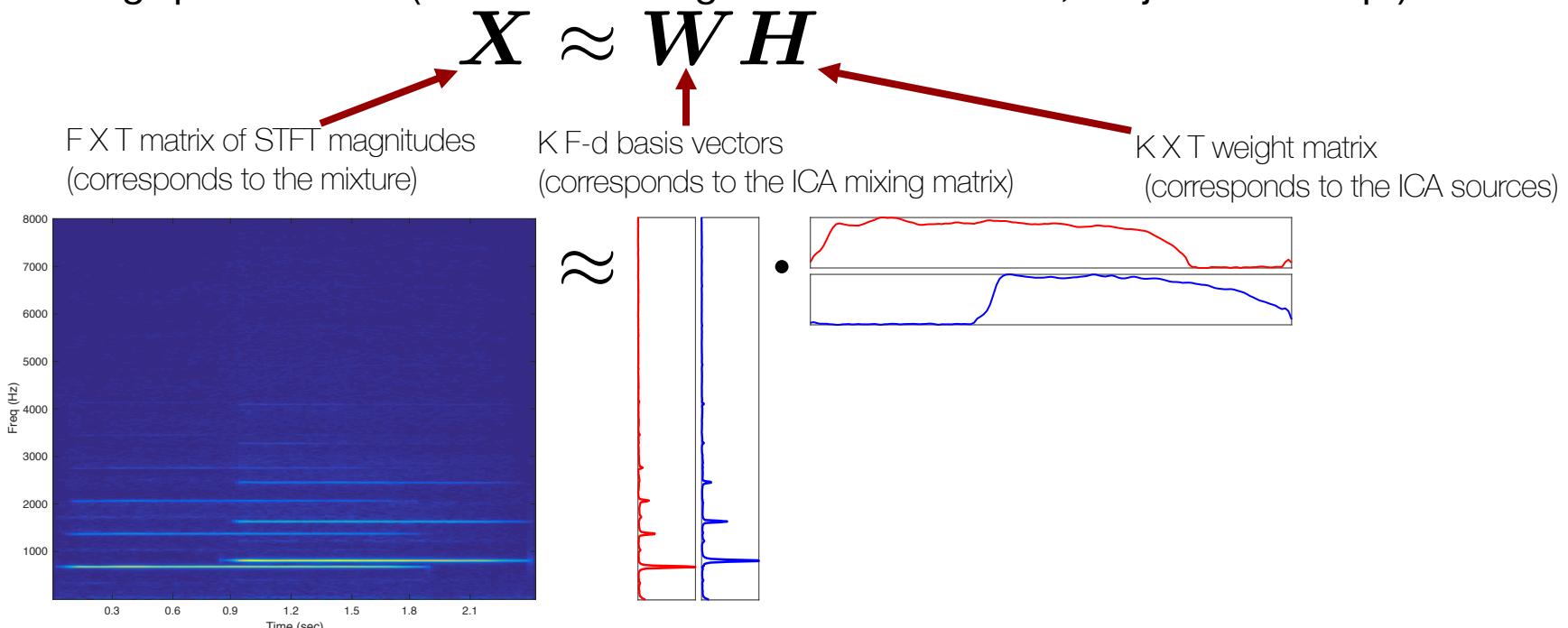


INDIANA UNIVERSITY

SCHOOL OF INFORMATICS, COMPUTING, AND ENGINEERING

Underdetermined Separation

- What if we have only one mixture
 - Can we still separate? How? (Hint: previous slide)
 - Increase the dimension via STFT on the 1D signal
 - And take the magnitudes
 - Then, the “mixing” procedure is (we’re not mixing the actual sources, it’s just a concept)



INDIANA UNIVERSITY

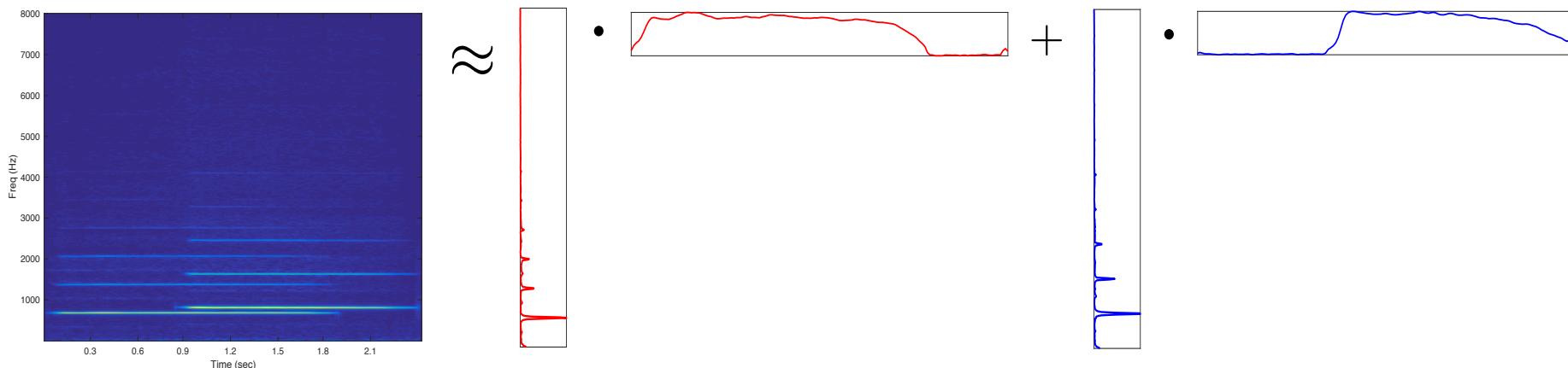
SCHOOL OF INFORMATICS, COMPUTING, AND ENGINEERING

Underdetermined Separation

- What if we have only one mixture
 - How can this procedure be source separation?
 - We can decompose the matrix into two components

$$X \approx WH$$

$$X \approx [w_1, w_2] \begin{bmatrix} h_1 \\ h_2 \end{bmatrix} = w_1 h_1 + w_2 h_2$$



- How do we find the best decomposition?
 - Well, PCA and ICA can actually do it to some degree

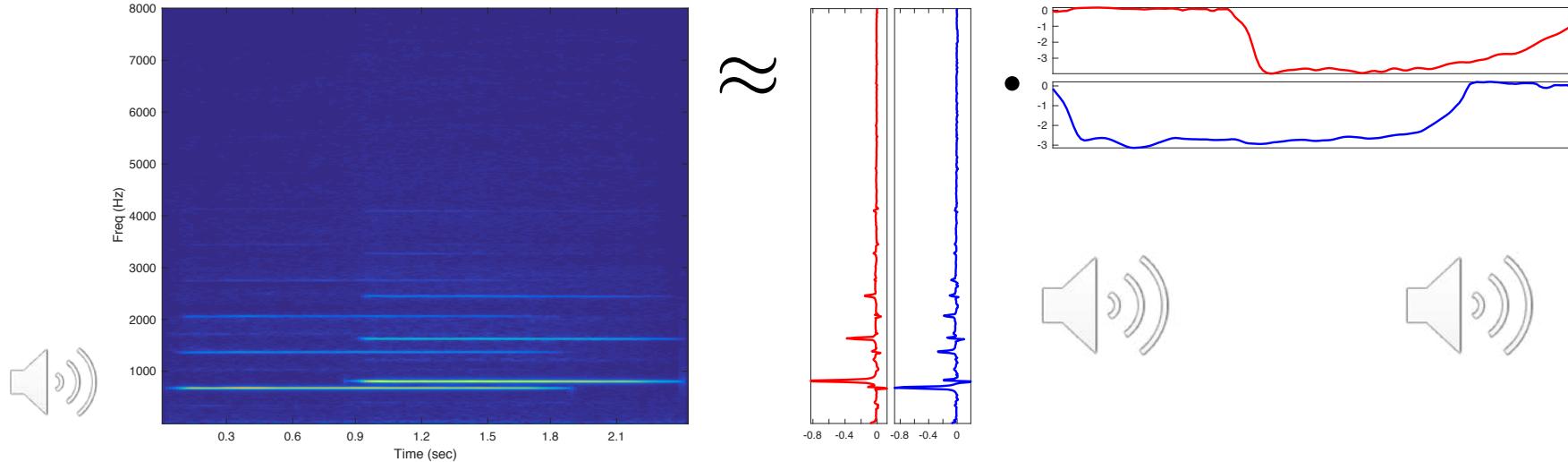


INDIANA UNIVERSITY

SCHOOL OF INFORMATICS, COMPUTING, AND ENGINEERING

Underdetermined Separation

- What if we have only one mixture
- The PCA results on the spectrogram of one mixture of two notes



- This seems to work, but..
 - It's a corner case where the source spectra are almost orthogonal (their inner product is near zero)
 - There's no guarantee
 - I don't like all the combination of positive and negative values to recover positive values

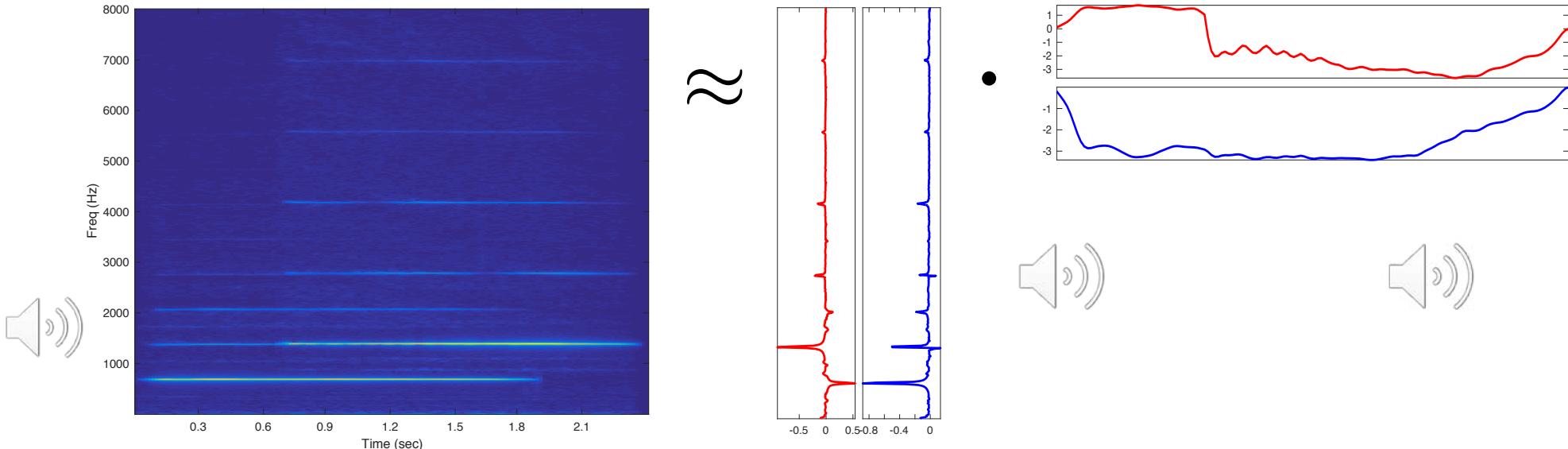


INDIANA UNIVERSITY

SCHOOL OF INFORMATICS, COMPUTING, AND ENGINEERING

Underdetermined Separation

- What if we have only one mixture
- Well, let's see what happens if the spectra are not orthogonal



- There are a lot of overlapping harmonics
 - Furthermore, there's no guarantee that each reconstruction is nonnegative

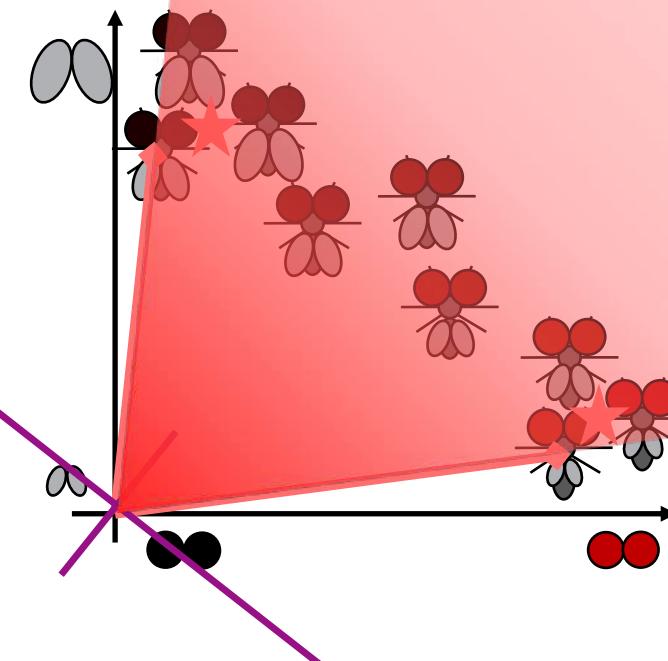


INDIANA UNIVERSITY

SCHOOL OF INFORMATICS, COMPUTING, AND ENGINEERING

Nonnegative Matrix Factorization

- Parts-based representation
 - The mating flies example once again
 - Originally it was a clustering problem
 - After mating
 - PCA can find a new subspace
 - With NMF we will find yet another one
 - NMF can find basis vectors that are
 - Nonnegative
 - So are their weights
 - Not required to be orthogonal
 - The convex cone defined by the basis is the subspace



INDIANA UNIVERSITY

SCHOOL OF INFORMATICS, COMPUTING, AND ENGINEERING

Nonnegative Matrix Factorization

- Parts-based representation

○ Objective function $\arg \min_{\mathbf{W}, \mathbf{H} \geq 0} \mathcal{J} = \frac{1}{2} \|\mathbf{X} - \mathbf{WH}\|_F^2 = \frac{1}{2} \text{tr} (\mathbf{X} - \mathbf{WH})^\top (\mathbf{X} - \mathbf{WH})$

○ Derivatives $\frac{\partial \mathcal{J}}{\partial \mathbf{W}} = \mathbf{WHH}^\top - \mathbf{XH}^\top$

$$\frac{\partial \mathcal{J}}{\partial \mathbf{H}} = \mathbf{W}^\top \mathbf{WH} - \mathbf{W}^\top \mathbf{X}$$

○ Gradients descent

$$\begin{aligned} \mathbf{W} &\leftarrow \mathbf{W} - \gamma(\mathbf{WHH}^\top - \mathbf{XH}^\top) \\ &= \mathbf{W} - \frac{\mathbf{W}}{\mathbf{WHH}^\top} \odot (\mathbf{WHH}^\top - \mathbf{XH}^\top) \quad \leftarrow \text{Choose the step size wisely} \end{aligned}$$

○ Multiplicative update rules

$$\mathbf{W} \leftarrow \mathbf{W} \odot \frac{\mathbf{XH}^\top}{\mathbf{WHH}^\top}$$

$$\mathbf{H} \leftarrow \mathbf{H} \odot \frac{\mathbf{W}^\top \mathbf{X}}{\mathbf{W}^\top \mathbf{WH}}$$

$$\begin{aligned} \frac{\partial \mathcal{J}}{\partial \mathbf{W}} &= \left[\frac{\partial \mathcal{J}}{\partial \mathbf{W}} \right]^+ - \left[\frac{\partial \mathcal{J}}{\partial \mathbf{W}} \right]^- \\ &= \mathbf{WHH}^\top - \mathbf{XH}^\top \end{aligned}$$

\leftarrow An alternative way to derive update rules



INDIANA UNIVERSITY

SCHOOL OF INFORMATICS, COMPUTING, AND ENGINEERING

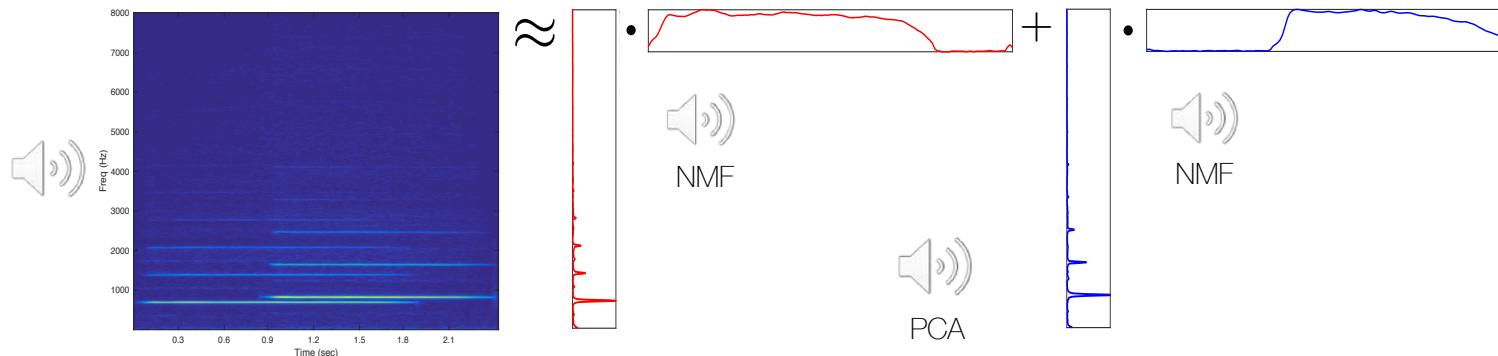
Lee&Seung Nature 1999, Lee&Seung NIPS 2001

Nonnegative Matrix Factorization

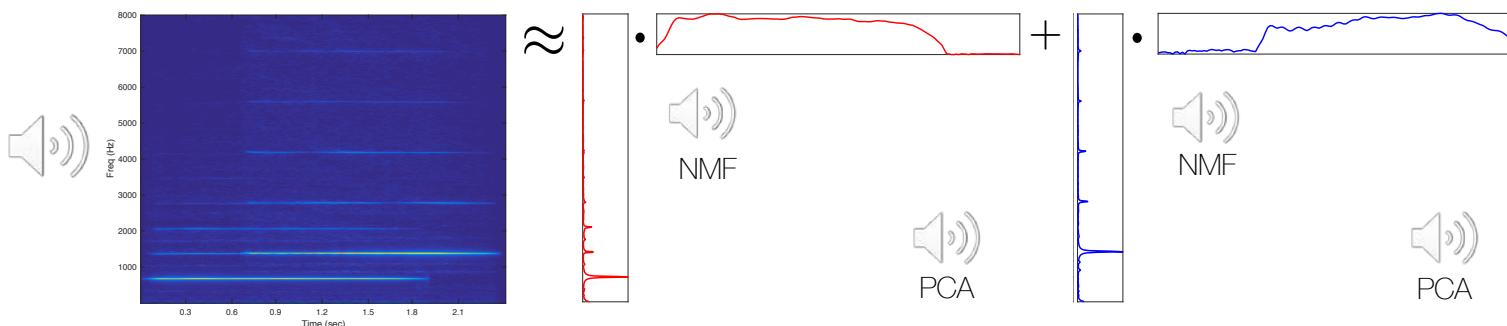
- Parts-based representation

- So, does it work?

- Below are from actual NMF runs
 - The almost orthogonal case



- The non-orthogonal case

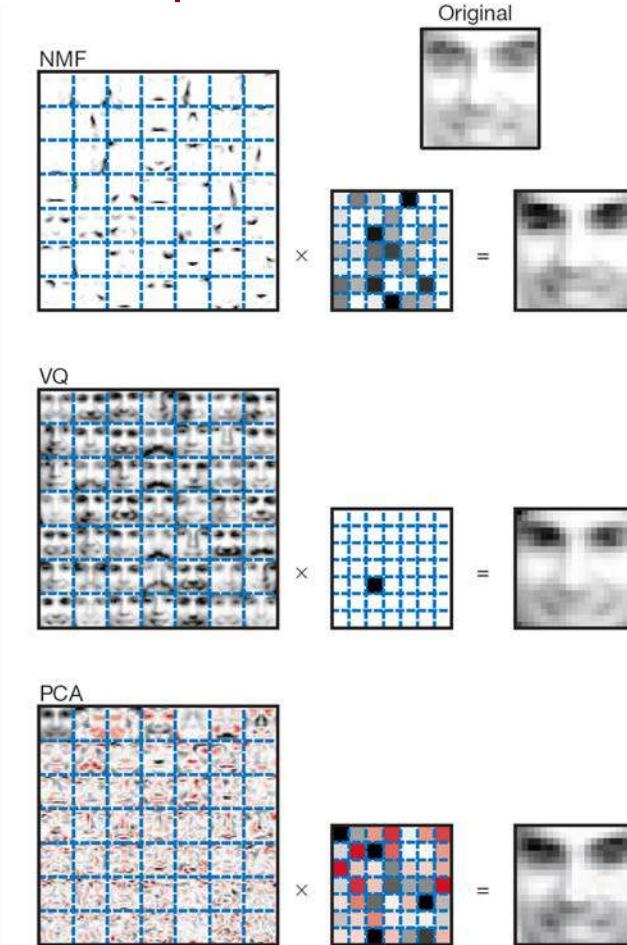


INDIANA UNIVERSITY

SCHOOL OF INFORMATICS, COMPUTING, AND ENGINEERING

Nonnegative Matrix Factorization

- Parts-based representation



NMF estimates parts-based representations,
something like Lego blocks.

Reconstruction is a linear combination of them,
but subtraction is not allowed.

VQ finds a bunch of cluster means.
Reconstruction is choosing the most similar mean.

PCA finds the holistic eigenfaces.
From the important one down to the subtle ones.



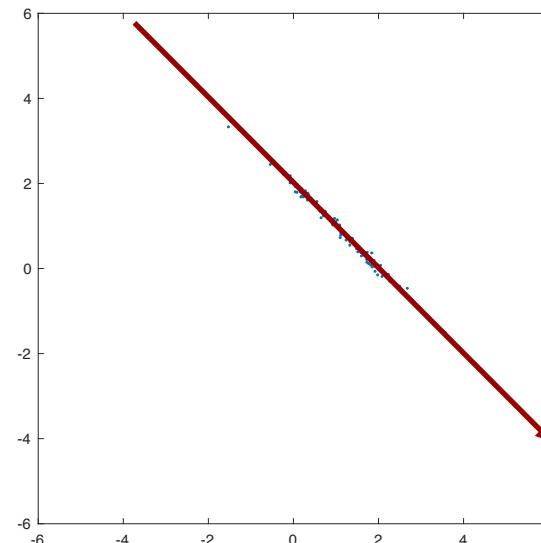
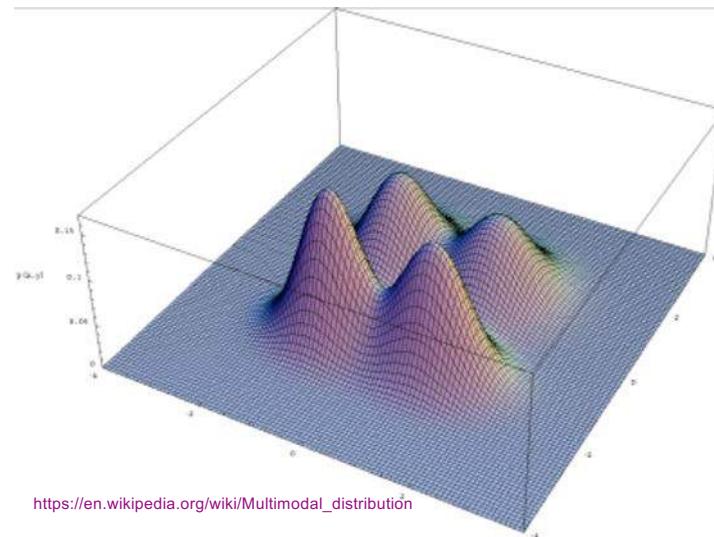
INDIANA UNIVERSITY

SCHOOL OF INFORMATICS, COMPUTING, AND ENGINEERING

Lee&Seung Nature 1999, Lee&Seung NIPS 2001

Recap

- Mixture models (clustering)
 - Assume less modes than the samples
 - Each sample belongs to one cluster, though there's uncertainty
 - Compression is done by keeping the cluster means as representatives
- Dimension reduction
 - Assumes the samples lie on a lower dimensional subspace
 - The subspace is defined with K basis vectors (so called dictionary)
 - $K < D$ (not always)
 - Various kinds of basis vectors depending on the constraints
 - Orthogonality, variance maximization, independence, nonnegativity, etc.



INDIANA UNIVERSITY

SCHOOL OF INFORMATICS, COMPUTING, AND ENGINEERING

Reading

- 6.1-6.7



INDIANA UNIVERSITY

SCHOOL OF INFORMATICS, COMPUTING, AND ENGINEERING



Thank You!



INDIANA UNIVERSITY

SCHOOL OF INFORMATICS, COMPUTING, AND ENGINEERING