

Psych 253

Advanced Statistical Modeling

Dimension Reduction

Daniel Yamins

Wu Tsai Neurosciences Institute
Departments of Psychology and Computer Science
Stanford Artificial Intelligence Laboratory
Stanford University

Russ Poldrack

Department of Psychology
Stanford University

Unsupervised data analysis

No “labels”, just “data”.

<http://scikit-learn.org/stable/modules/classes.html#module-sklearn.decomposition>

Dimension reduction: `sklearn.decomposition`

Unsupervised data analysis

No “labels”, just “data”.

<http://scikit-learn.org/stable/modules/classes.html#module-sklearn.decomposition>

Dimension reduction: **sklearn.decomposition**

Principal Components Analysis (PCA)

Factor Analysis

Multi-Dimensional Scaling (MDS)

Independent Components Analysis (ICA)

Non-negative Matrix Factorization (NMF)

Unsupervised data analysis

No “labels”, just “data”.

<http://scikit-learn.org/stable/modules/classes.html#module-sklearn.decomposition>

Dimension reduction: **sklearn.decomposition**

Principal Components Analysis (PCA)

Factor Analysis

Multi-Dimensional Scaling (MDS)

Independent Components Analysis (ICA)

Non-negative Matrix Factorization (NMF)

Unsupervised data analysis

No “labels”, just “data”.

<http://scikit-learn.org/stable/modules/classes.html#module-sklearn.decomposition>

Dimension reduction: **sklearn.decomposition**

Principal Components Analysis (PCA)

Factor Analysis

Multi-Dimensional Scaling (MDS)

Independent Components Analysis (ICA)

Non-negative Matrix Factorization (NMF)

<http://scikit-learn.org/stable/modules/classes.html#module-sklearn.cluster>

Clustering: **sklearn.cluster**

Unsupervised data analysis

No “labels”, just “data”.

<http://scikit-learn.org/stable/modules/classes.html#module-sklearn.decomposition>

Dimension reduction: **sklearn.decomposition**

Principal Components Analysis (PCA)

Factor Analysis

Multi-Dimensional Scaling (MDS)

Independent Components Analysis (ICA)

Non-negative Matrix Factorization (NMF)

<http://scikit-learn.org/stable/modules/classes.html#module-sklearn.cluster>

Clustering: **sklearn.cluster**

k-means clustering

Hierarchical k-means

Affinity Propagation

DBScan

Unsupervised data analysis

No “labels”, just “data”.

<http://scikit-learn.org/stable/modules/classes.html#module-sklearn.decomposition>

Dimension reduction: **sklearn.decomposition**

Principal Components Analysis (PCA)

Factor Analysis

Multi-Dimensional Scaling (MDS)

Independent Components Analysis (ICA)

Non-negative Matrix Factorization (NMF)

<http://scikit-learn.org/stable/modules/classes.html#module-sklearn.cluster>

Clustering: **sklearn.cluster**

k-means clustering

Hierarchical k-means

Affinity Propagation

DBScan

Last time

Unsupervised data analysis

No “labels”, just “data”.

<http://scikit-learn.org/stable/modules/classes.html#module-sklearn.decomposition>

Dimension reduction: **sklearn.decomposition**

Principal Components Analysis (PCA — and CCA)

Factor Analysis

Multi-Dimensional Scaling (MDS)

Independent Components Analysis (ICA)

Non-negative Matrix Factorization (NMF)

Today

<http://scikit-learn.org/stable/modules/classes.html#module-sklearn.cluster>

Clustering: **sklearn.cluster**

k-means clustering

Hierarchical k-means

Affinity Propagation

DBScan

Principal Components Analysis

1) Data matrix **X** of shape (m, n)

Principal Components Analysis

- 1) Data matrix **X** of shape (m, n) — mean centered!
- 2) Covariance matrix **$X^T X$** of shape

Principal Components Analysis

- 1) Data matrix \mathbf{X} of shape (m, n) — mean centered!
- 2) Covariance matrix $\mathbf{X}^T\mathbf{X}$ of shape (n, n)

Principal Components Analysis

- 1) Data matrix \mathbf{X} of shape (m, n) — mean centered!
- 2) Covariance matrix $\mathbf{X}^T\mathbf{X}$ of shape (n, n)
- 3) Compute eigenvalues $\mathbf{d}_1 \mathbf{d}_2 \dots \mathbf{d}_n$ and eigenvectors $\mathbf{w}_1 \mathbf{w}_2 \dots \mathbf{w}_n$ of $\mathbf{X}^T\mathbf{X}$

ordered from biggest to smallest

Principal Components Analysis

- 1) Data matrix \mathbf{X} of shape (m, n) — mean centered!
 - 2) Covariance matrix $\mathbf{X}^T\mathbf{X}$ of shape (n, n)
 - 3) Compute eigenvalues $\mathbf{d}_1 \mathbf{d}_2 \dots \mathbf{d}_n$ and eigenvectors $\mathbf{w}_1 \mathbf{w}_2 \dots \mathbf{w}_n$ of $\mathbf{X}^T\mathbf{X}$
- ↑
ordered from biggest to smallest

\mathbf{w}_1 = principal axis

$\mathbf{X} \cdot \vec{w}_1$ = projection of data onto principal axis

Principal Components Analysis

- 1) Data matrix \mathbf{X} of shape (m, n) — mean centered!
- 2) Covariance matrix $\mathbf{X}^T\mathbf{X}$ of shape (n, n)
- 3) Compute eigenvalues $\mathbf{d}_1 \mathbf{d}_2 \dots \mathbf{d}_n$ and eigenvectors $\mathbf{w}_1 \mathbf{w}_2 \dots \mathbf{w}_n$ of $\mathbf{X}^T\mathbf{X}$



ordered from biggest to smallest

\mathbf{w}_1 = principal axis $\mathbf{X} \cdot \vec{w}_1$ = projection of data onto principal axis

\mathbf{w}_2 = next principal axis

\mathbf{w}_k = k-th principal axis

Principal Components Analysis

- 1) Data matrix \mathbf{X} of shape (m, n) — mean centered!
- 2) Covariance matrix $\mathbf{X}^T\mathbf{X}$ of shape (n, n)
- 3) Compute eigenvalues $\mathbf{d}_1 \mathbf{d}_2 \dots \mathbf{d}_n$ and eigenvectors $\mathbf{w}_1 \mathbf{w}_2 \dots \mathbf{w}_n$ of $\mathbf{X}^T\mathbf{X}$

\uparrow
 ordered from biggest to smallest

\mathbf{w}_1 = principal axis $\mathbf{X} \cdot \vec{w}_1$ = projection of data onto principal axis

\mathbf{w}_2 = next principal axis $\mathbf{X} \cdot \vec{w}_2$ = proj. of data onto next principal axis

\mathbf{w}_k = k-th principal axis

Principal Components Analysis

- 1) Data matrix \mathbf{X} of shape (m, n) — mean centered!
- 2) Covariance matrix $\mathbf{X}^T\mathbf{X}$ of shape (n, n)
- 3) Compute eigenvalues $\mathbf{d}_1 \mathbf{d}_2 \dots \mathbf{d}_n$ and eigenvectors $\mathbf{w}_1 \mathbf{w}_2 \dots \mathbf{w}_n$ of $\mathbf{X}^T\mathbf{X}$

\uparrow
 ordered from biggest to smallest

\mathbf{w}_1 = principal axis $\mathbf{X} \cdot \vec{w}_1$ = projection of data onto principal axis

\mathbf{w}_2 = next principal axis $\mathbf{X} \cdot \vec{w}_2$ = proj. of data onto next principal axis

\mathbf{w}_k = k-th principal axis $\mathbf{X} \cdot [\vec{w}_1, \dots, \vec{w}_k]$

Principal Components Analysis

- 1) Data matrix \mathbf{X} of shape (m, n) — mean centered!
- 2) Covariance matrix $\mathbf{X}^T\mathbf{X}$ of shape (n, n)
- 3) Compute eigenvalues $\mathbf{d}_1 \mathbf{d}_2 \dots \mathbf{d}_n$ and eigenvectors $\mathbf{w}_1 \mathbf{w}_2 \dots \mathbf{w}_n$ of $\mathbf{X}^T\mathbf{X}$

\uparrow
 ordered from biggest to smallest

\mathbf{w}_1 = principal axis $\mathbf{X} \cdot \vec{w}_1$ = projection of data onto principal axis

\mathbf{w}_2 = next principal axis $\mathbf{X} \cdot \vec{w}_2$ = proj. of data onto next principal axis

\mathbf{w}_k = k-th principal axis $\mathbf{X} \cdot [\vec{w}_1, \dots, \vec{w}_k]$
 = projection of data onto top \mathbf{k} principal axes

Principal Components Analysis

- 1) Data matrix \mathbf{X} of shape (m, n) — mean centered!
- 2) Covariance matrix $\mathbf{X}^T\mathbf{X}$ of shape (n, n)
- 3) Compute eigenvalues $\mathbf{d}_1 \mathbf{d}_2 \dots \mathbf{d}_n$ and eigenvectors $\mathbf{w}_1 \mathbf{w}_2 \dots \mathbf{w}_n$ of $\mathbf{X}^T\mathbf{X}$

 ordered from biggest to smallest

\mathbf{w}_1 = principal axis $\mathbf{X} \cdot \vec{w}_1$ = projection of data onto principal axis

\mathbf{w}_2 = next principal axis $\mathbf{X} \cdot \vec{w}_2$ = proj. of data onto next principal axis

\mathbf{w}_k = k-th principal axis $\mathbf{X} \cdot [\vec{w}_1, \dots, \vec{w}_k]$  $\frac{\sum_{i \leq k} d_i}{\sum_{i=1}^n d_i}$
 fraction of variance of original data explained
 = projection of data onto top \mathbf{k} principal axes

[Ipython: principal components analysis]

Normal projection of PCA is:

$$[p_1, \dots, p_k] = X \cdot [\vec{w}_1, \dots, \vec{w}_k]$$

Normal projection of PCA is:

$$[p_1, \dots, p_k] = X \cdot [\vec{w}_1, \dots, \vec{w}_k]$$

The “whitened” projection of the PCA is:

$$[p_1^{wh}, \dots, p_k^{wh}] = (X \cdot [\vec{w}_1, \dots, \vec{w}_k]) / [\sqrt{d_1}, \dots, \sqrt{d_k}]$$

Normal projection of PCA is:

$$[p_1, \dots, p_k] = X \cdot [\vec{w}_1, \dots, \vec{w}_k]$$

The “whitened” projection of the PCA is:

$$[p_1^{wh}, \dots, p_k^{wh}] = (X \cdot [\vec{w}_1, \dots, \vec{w}_k]) / [\sqrt{d_1}, \dots, \sqrt{d_k}]$$

What does this accomplish?

Normal projection of PCA is:

$$[p_1, \dots, p_k] = X \cdot [\vec{w}_1, \dots, \vec{w}_k]$$

The “whitened” projection of the PCA is:

$$[p_1^{wh}, \dots, p_k^{wh}] = (X \cdot [\vec{w}_1, \dots, \vec{w}_k]) / [\sqrt{d_1}, \dots, \sqrt{d_k}]$$

What does this accomplish? Makes the amount of variance = 1 in each dimension

PCA

Normal projection of PCA is:

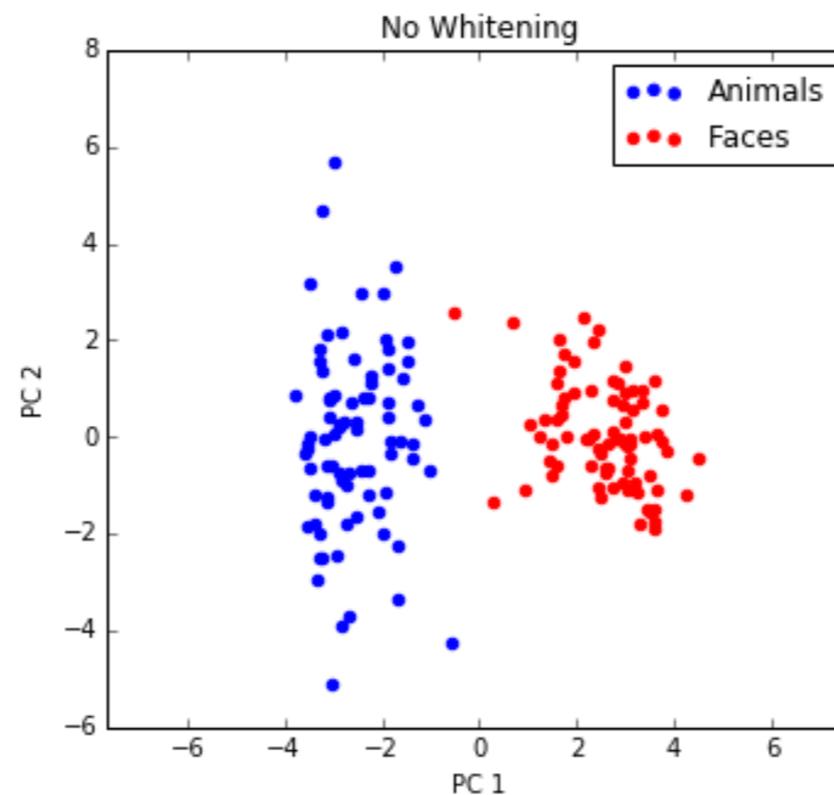
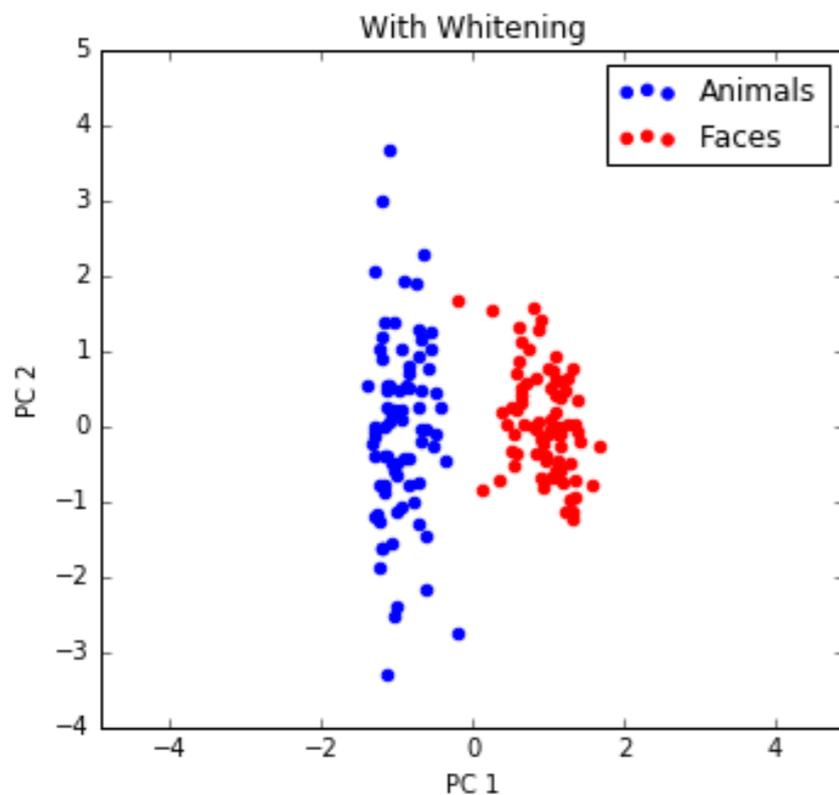
$$[p_1, \dots, p_k] = X \cdot [\vec{w}_1, \dots, \vec{w}_k]$$

The “whitened” projection of the PCA is:

$$[p_1^{wh}, \dots, p_k^{wh}] = (X \cdot [\vec{w}_1, \dots, \vec{w}_k]) / [\sqrt{d_1}, \dots, \sqrt{d_k}]$$

What does this accomplish?

Makes the amount of variance = 1 in each dimension



PCA

Normal projection of PCA is:

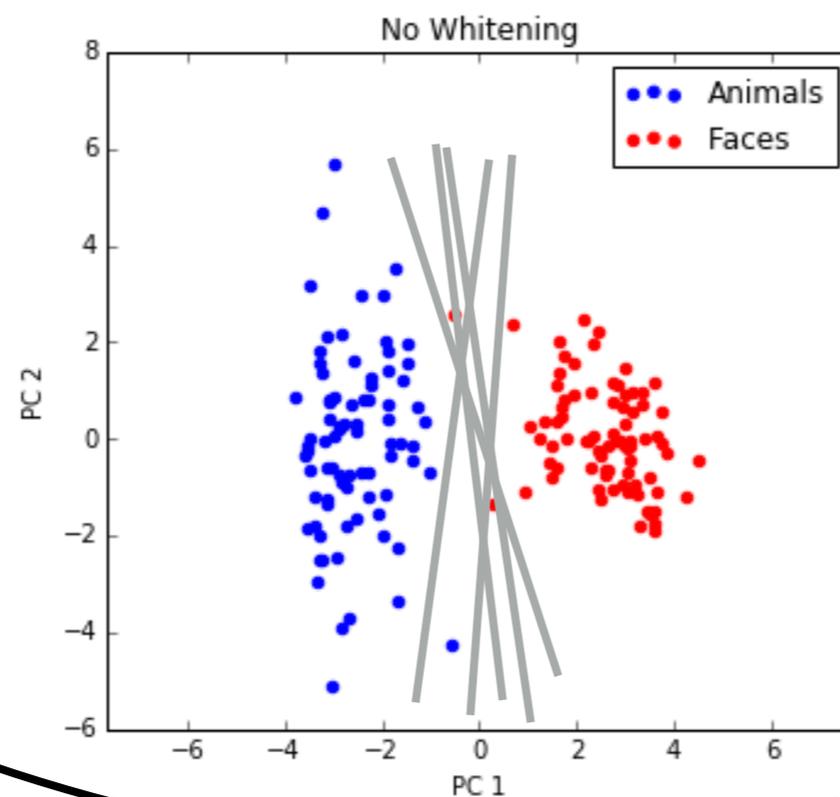
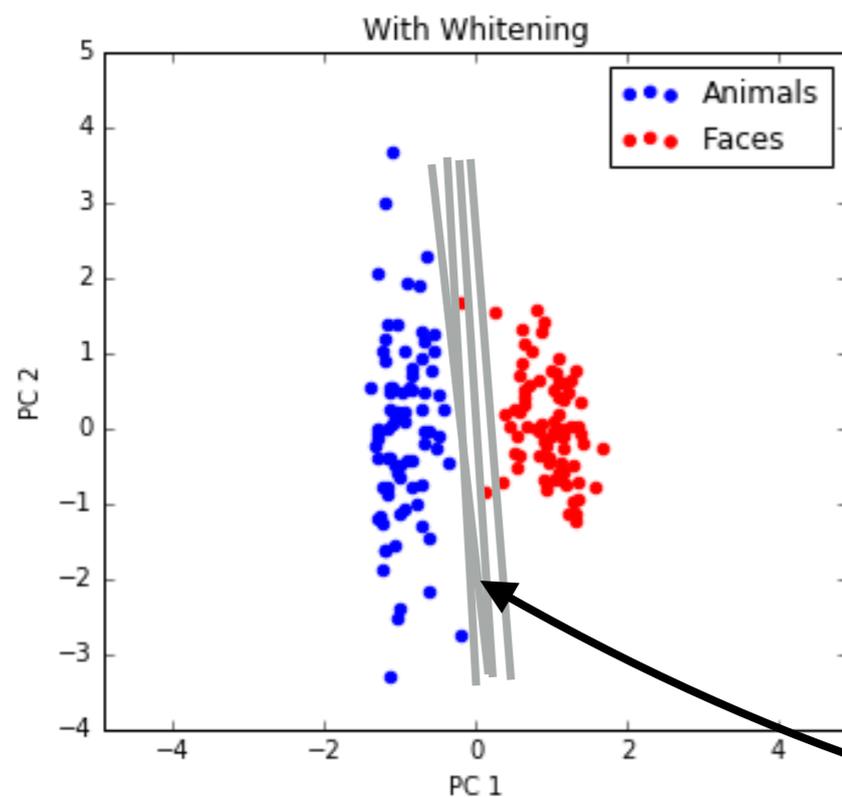
$$[p_1, \dots, p_k] = X \cdot [\vec{w}_1, \dots, \vec{w}_k]$$

The “whitened” projection of the PCA is:

$$[p_1^{wh}, \dots, p_k^{wh}] = (X \cdot [\vec{w}_1, \dots, \vec{w}_k]) / [\sqrt{d_1}, \dots, \sqrt{d_k}]$$

What does this accomplish?

Makes the amount of variance = 1 in each dimension



Whitening can make things in downstream operations better.

Such as regularization — less mischief is possible.

Recover Principal Axes of Data

Basic noisy compressive model:

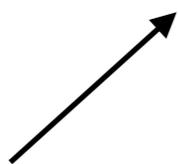
$$X = W \cdot P + \mu + \epsilon$$

Recover Principal Axes of Data

Basic noisy compressive model:

$$X = W \cdot P + \mu + \epsilon$$

data vec
of shape
n



Recover Principal Axes of Data

Basic noisy compressive model:

$$X = W \cdot P + \mu + \epsilon$$

data vec
of shape
n

k vectors
of shape
(n, k)

The diagram shows the equation $X = W \cdot P + \mu + \epsilon$. Two arrows originate from the text below. The first arrow points from the text 'data vec of shape n' to the variable X . The second arrow points from the text 'k vectors of shape (n, k)' to the variable W .

Recover Principal Axes of Data

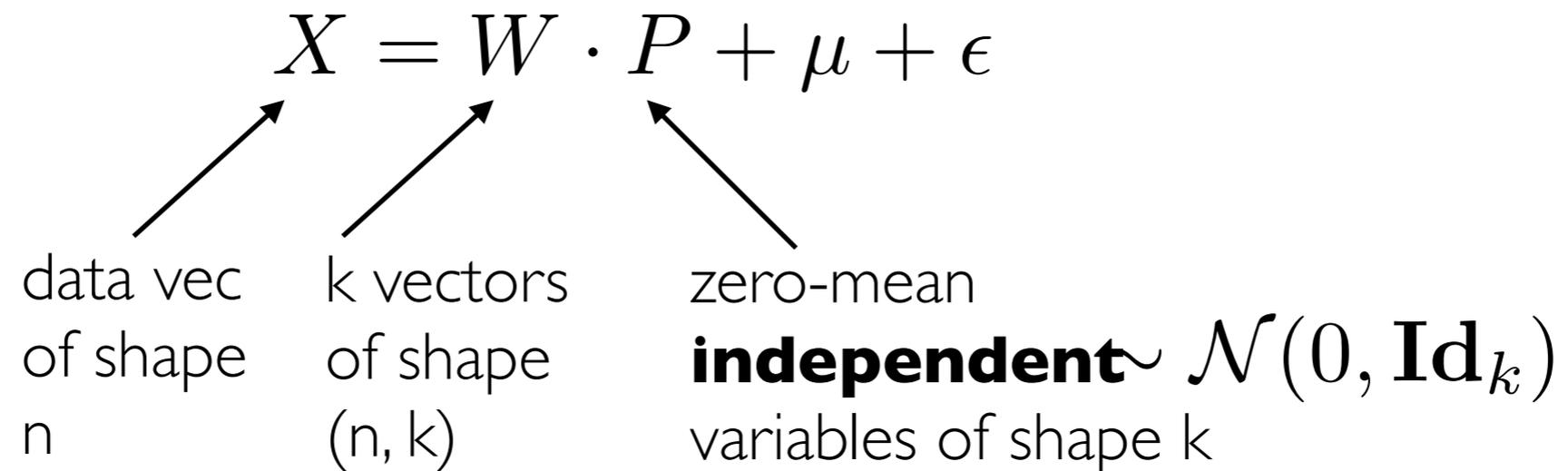
Basic noisy compressive model:

$$X = W \cdot P + \mu + \epsilon$$

data vec
of shape
n

k vectors
of shape
(n, k)

zero-mean
independent $\sim \mathcal{N}(0, \mathbf{Id}_k)$
variables of shape k



Recover Principal Axes of Data

Basic noisy compressive model:

$$X = W \cdot P + \mu + \epsilon \sim \mathcal{N}(0, \sigma^2 \mathbf{Id}_n)$$

data vec of shape n

k vectors of shape (n, k)

vector of means of shape n

n **independent** gaussian random variables, all of same variance

zero-mean **independent** variables of shape k

$$\sim \mathcal{N}(0, \mathbf{Id}_k)$$

Recover Principal Axes of Data

Basic noisy compressive model:

$$X = W \cdot P + \mu + \epsilon$$

vector of means of shape n

data vec
of shape
n

k vectors
of shape
(n, k)

zero-mean
independent $\sim \mathcal{N}(0, \mathbf{Id}_k)$
variables of shape k

The diagram shows the equation $X = W \cdot P + \mu + \epsilon$ centered on the page. Four arrows point from descriptive text to the terms in the equation:
1. An arrow points from 'data vec of shape n' to the variable X .
2. An arrow points from 'k vectors of shape (n, k)' to the matrix W .
3. An arrow points from 'zero-mean independent $\sim \mathcal{N}(0, \mathbf{Id}_k)$ variables of shape k' to the variable P .
4. An arrow points from 'vector of means of shape n' to the vector μ .

Recover Principal Axes of Data

Basic noisy compressive model:

$$X = W \cdot P + \mu + \epsilon \sim \mathcal{N}(0, \sigma^2 \mathbf{Id}_n)$$

data vec of shape n

k vectors of shape (n, k)

vector of means of shape n

n **independent** gaussian random variables, all of same variance

zero-mean **independent** variables of shape k

$$\sim \mathcal{N}(0, \mathbf{Id}_k)$$

COV =

Recover Principal Axes of Data

Basic noisy compressive model:

$$X = W \cdot P + \mu + \epsilon \sim \mathcal{N}(0, \sigma^2 \mathbf{Id}_n)$$

vector of means of shape n

data vec of shape n

k vectors of shape (n, k)

zero-mean **independent** $\sim \mathcal{N}(0, \mathbf{Id}_k)$ variables of shape k

n **independent** gaussian random variables, all of same variance

$$\mathbf{cov} = (X - \mu)^T (X - \mu)$$

Recover Principal Axes of Data

Basic noisy compressive model:

$$X = W \cdot P + \mu + \epsilon \sim \mathcal{N}(0, \sigma^2 \mathbf{Id}_n)$$

vector of means of shape n

data vec of shape n

k vectors of shape (n, k)

zero-mean **independent** $\sim \mathcal{N}(0, \mathbf{Id}_k)$ variables of shape k

n **independent** gaussian random variables, all of same variance

$$\mathbf{cov} = (X - \mu)^T (X - \mu) = W^T W + \sigma^2 \mathbf{Id}_n$$

Recover Principal Axes of Data

Basic noisy compressive model:

$$X = W \cdot P + \mu + \epsilon \sim \mathcal{N}(0, \sigma^2 \mathbf{Id}_n)$$

vector of means of shape n

data vec of shape n

k vectors of shape (n, k)

zero-mean **independent** $\sim \mathcal{N}(0, \mathbf{Id}_k)$ variables of shape k

n **independent** gaussian random variables, all of same variance

$$\mathbf{cov} = (X - \mu)^T (X - \mu) = W^T W + \sigma^2 \mathbf{Id}_n$$

If you choose like this . . .

$$W = U(D_k - \sigma^2 \mathbf{Id}_k)^{1/2}$$

U = top k eigenvectors of **cov**

D_k = top k eigenvalues of **cov**

. . . then best approx. to X

Recover Principal Axes of Data

Basic noisy compressive model:

$$X = W \cdot P + \mu + \epsilon \sim \mathcal{N}(\mu, \sigma^2 \mathbf{Id}_n)$$

data vec of shape n \nearrow X \nwarrow vector of means of shape n μ
 k vectors of shape (n, k) \nearrow W \nwarrow 0-vector of length n ϵ
 $P \sim \mathcal{N}(0, \mathbf{Id}_k)$ zero-mean independent variables of shape k
 $\epsilon \sim \mathcal{N}(0, \sigma^2 \mathbf{Id}_n)$ n **independent** gaussian random variables, all of same variance

$$\mathbf{cov} = (X - \mu)^T (X - \mu) = W^T W + \sigma^2 \mathbf{Id}_n$$

If you choose like this . . .

$$W = U(D_k - \sigma^2 \mathbf{Id}_k)^{1/2}$$

$$\sigma^2 = \frac{1}{n - k} \sum_{i > k} d_i$$

remaining variance is treated as noise

U = top k eigenvectors of **cov**

D_k = top k eigenvalues of **cov**

. . . then best approx. to X

Factor Analysis

Factor Analysis is the *same* as PCA, except

$$X = W \cdot P + \mu + \epsilon \sim \mathcal{N}(0, \Psi_n)$$

vector of means of shape n

data vec of shape n

k vectors of shape (n, k)

zero-mean **independent** variables of shape k

$\sim \mathcal{N}(0, \mathbf{Id}_k)$

n independent gaussian random variables, **different** variances

Ψ_n = diagonal noise covariance matrix

Factor Analysis

Factor Analysis is the *same* as PCA, except

$$X = W \cdot P + \mu + \epsilon \sim \mathcal{N}(0, \Psi_n)$$

vector of means of shape n

data vec of shape n

k vectors of shape (n, k)

zero-mean **independent** variables of shape k

$\sim \mathcal{N}(0, \mathbf{Id}_k)$

n independent gaussian random variables, **different** variances

Ψ_n = diagonal noise covariance matrix

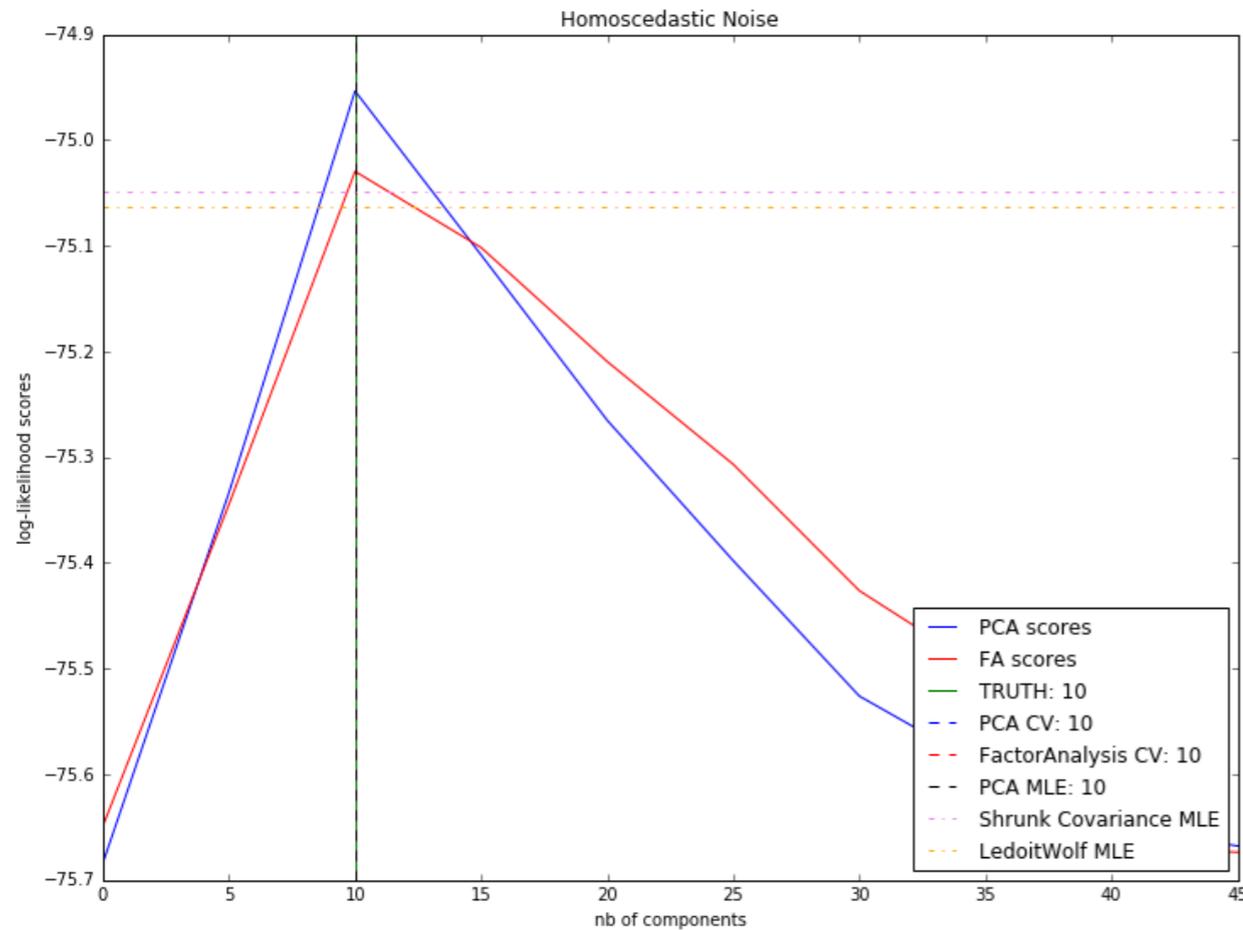
Use Factor Analysis if:

the noise across your prediction dimensions is not the same magnitude

Factor Analysis

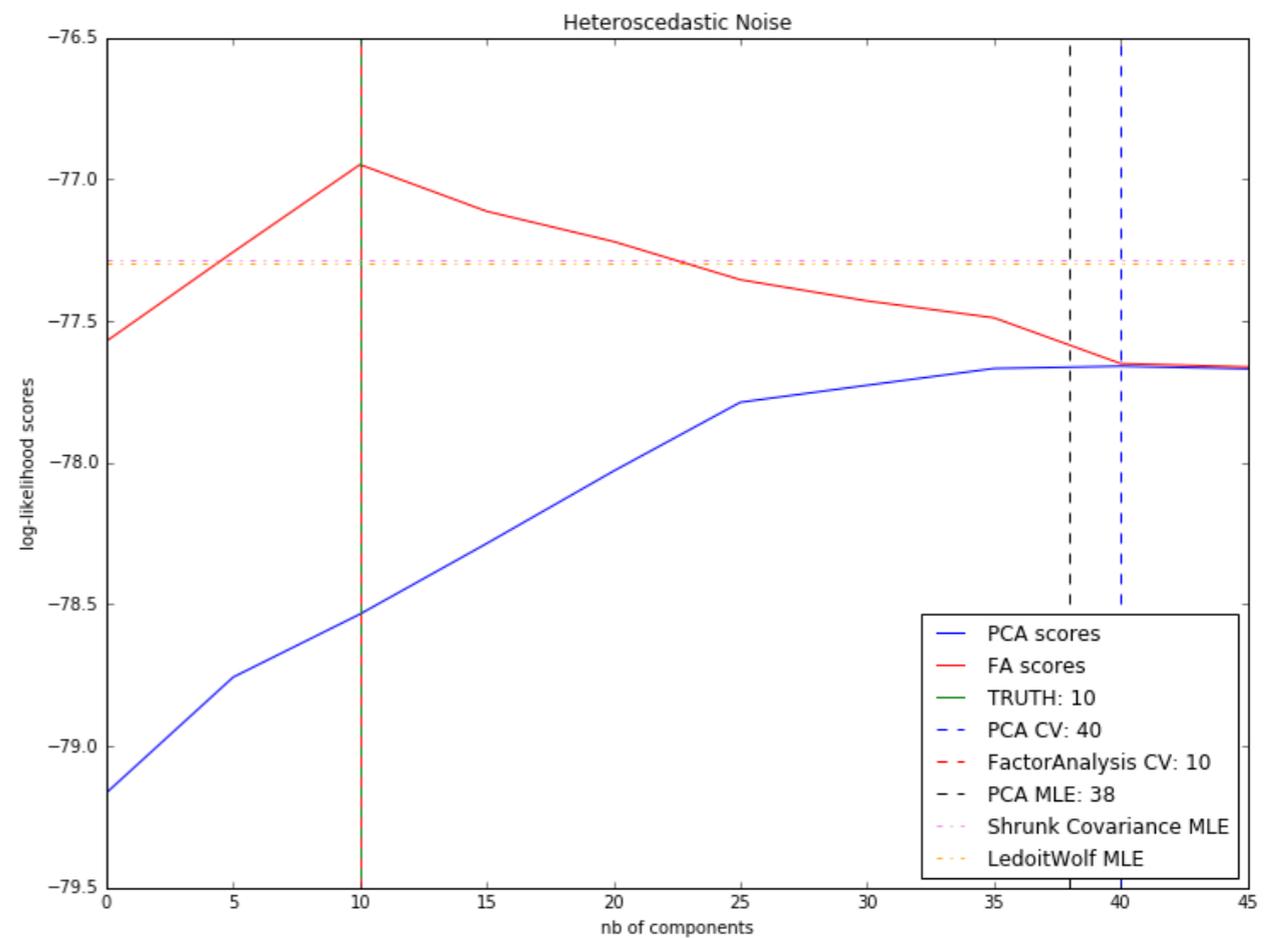
“Homoscedastic” noise:

$$\mathcal{N}(0, \sigma^2 \mathbf{Id}_n)$$



“Heteroscedastic” noise:

$$\mathcal{N}(0, \Psi_n)$$

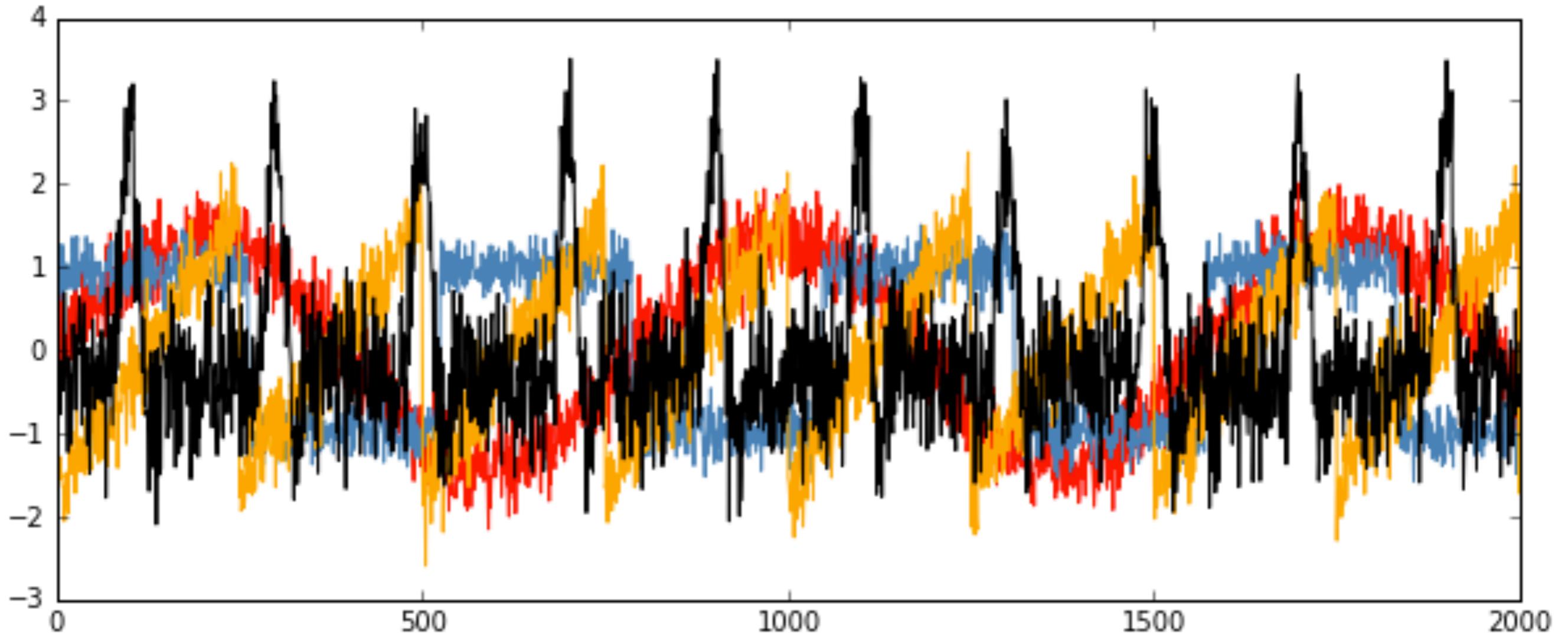


Use Factor Analysis if:

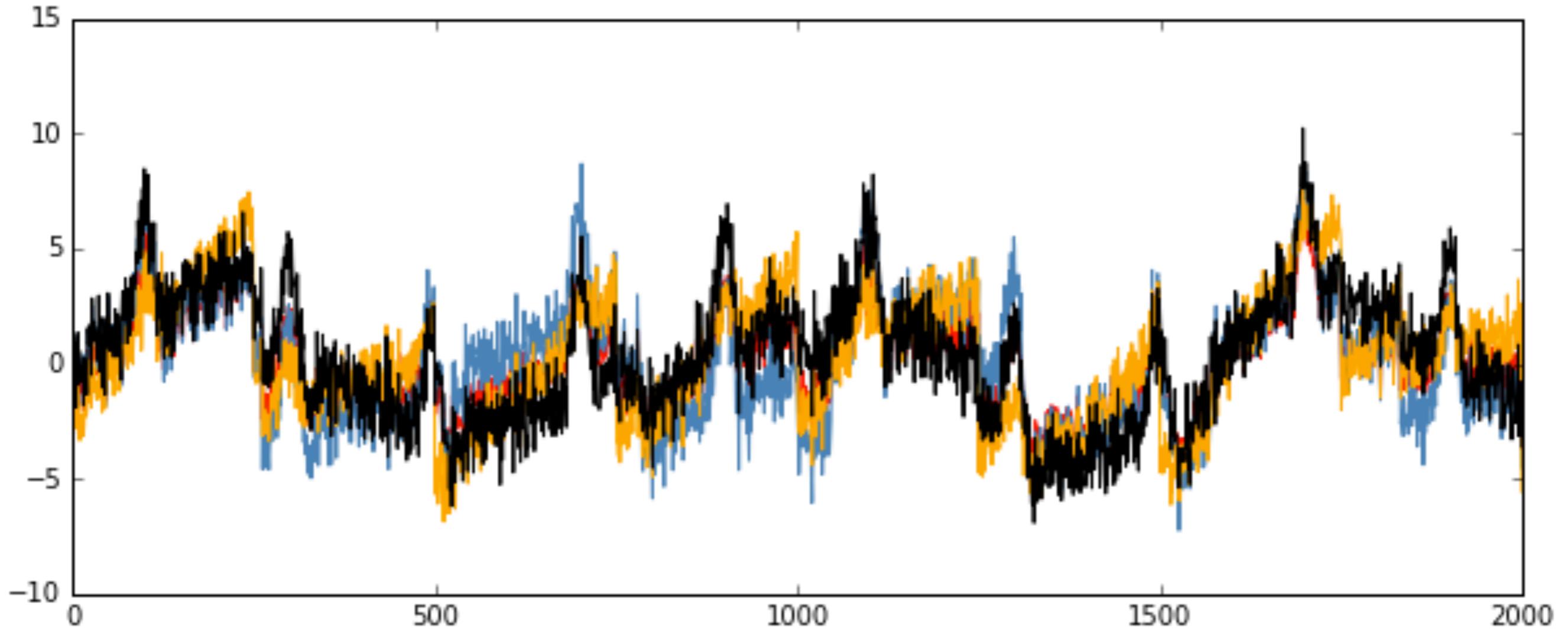
the noise across your prediction dimensions is not the same magnitude

(“heteroscedastic” noise)

Sources



Mixed Signals

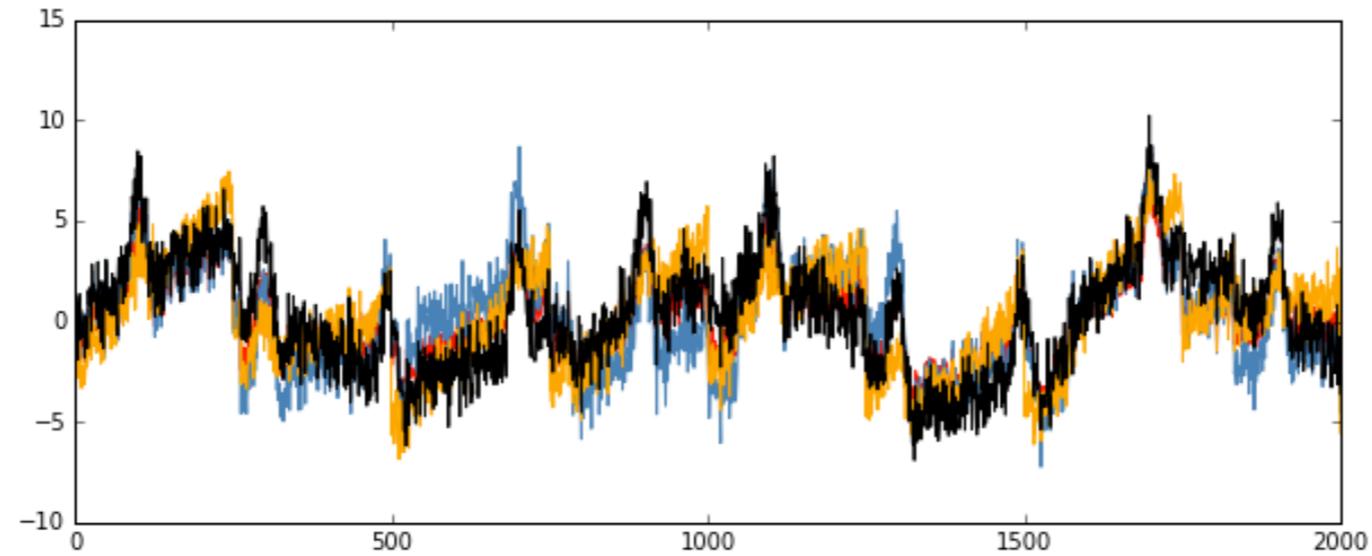


$$X = S \cdot A^T$$

Signals Sources Mixing Matrix

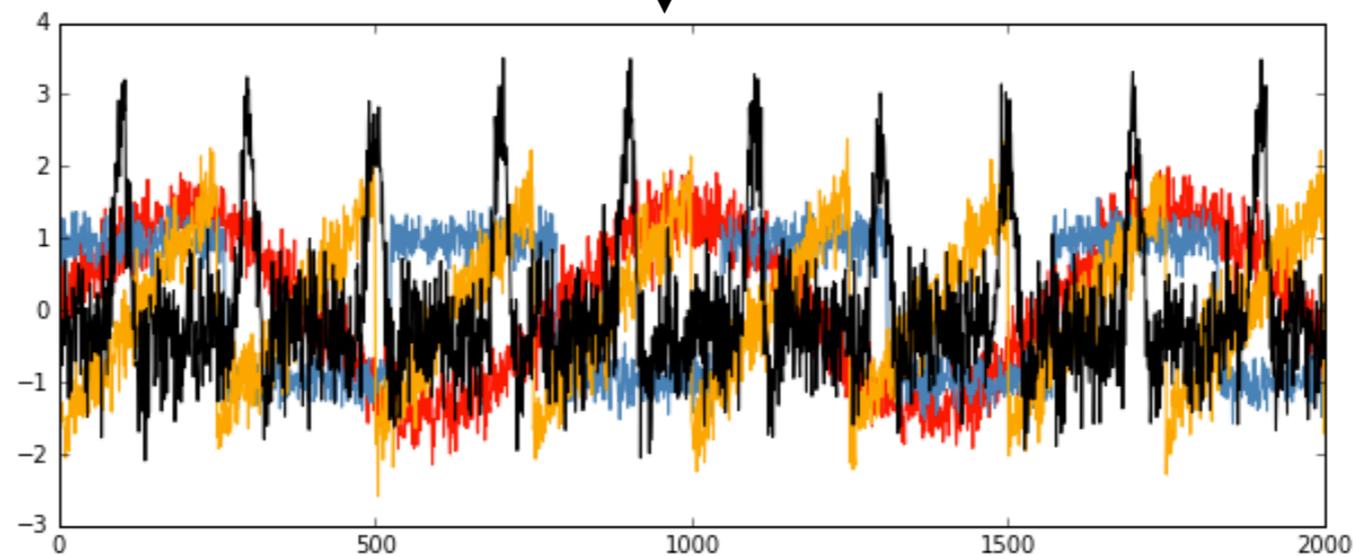
Independent component analysis

Signals



$$X = S \cdot A^T$$

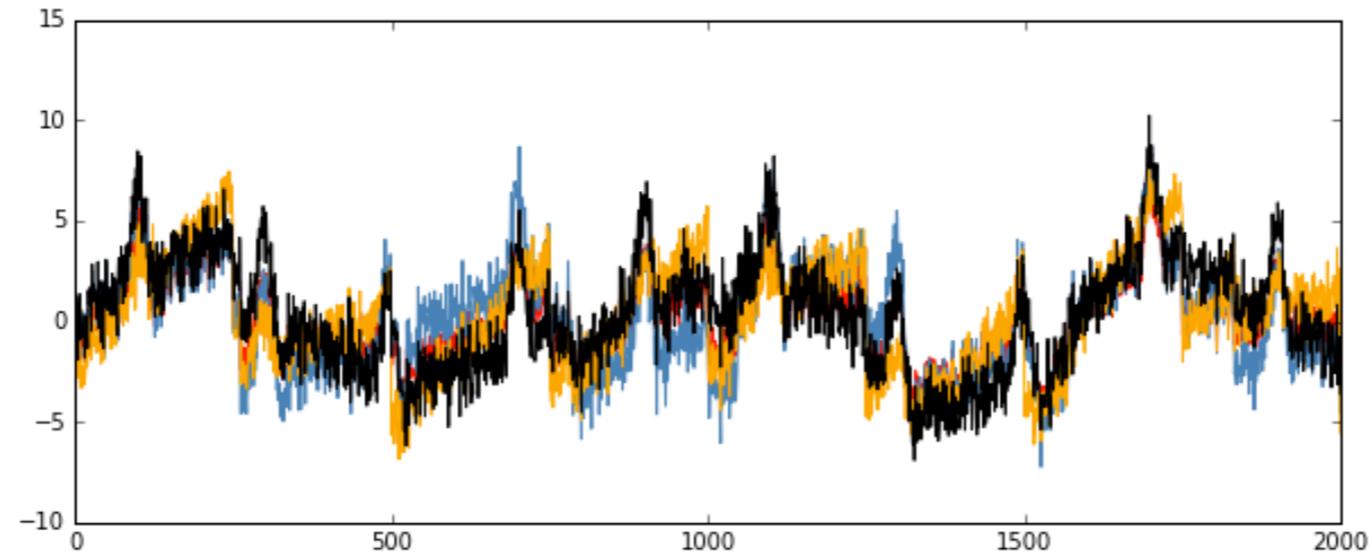
How to recover the sources??



Sources

Independent component analysis

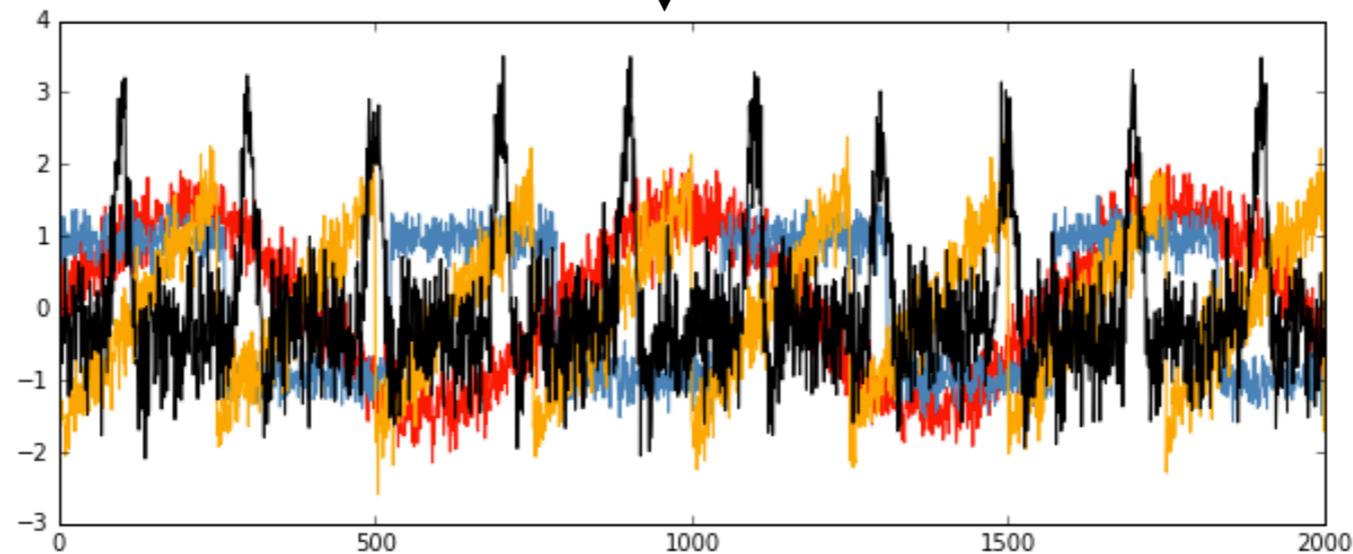
Signals



$$X = S \cdot A^T$$

How to recover the sources??

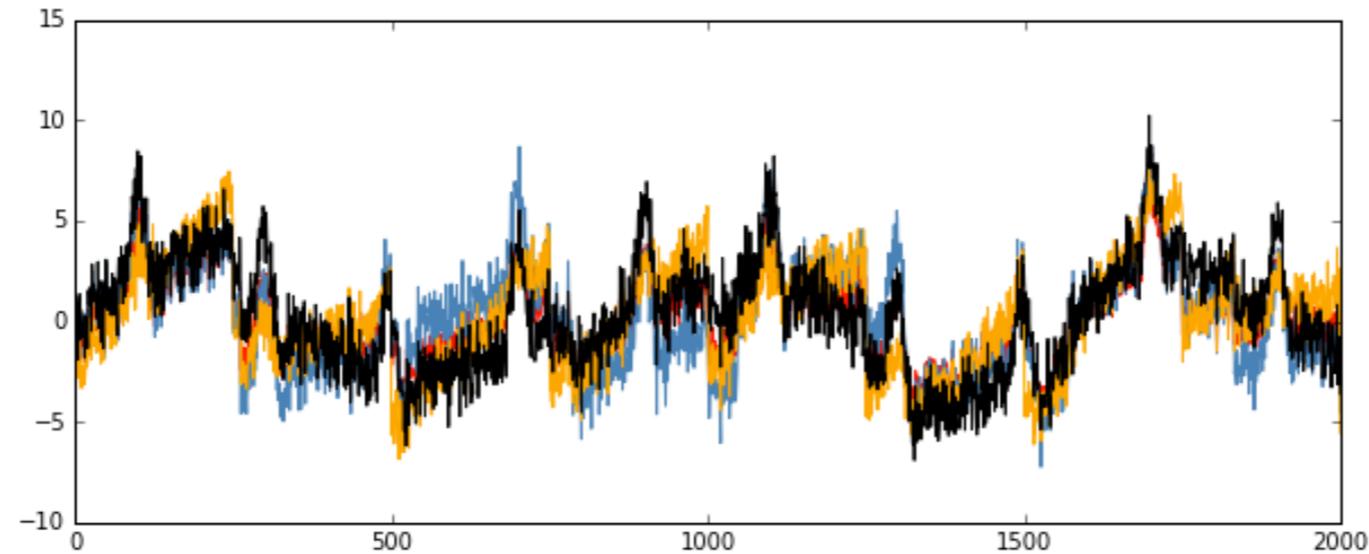
Problem: **S**, **A** unknown



Sources

Independent component analysis

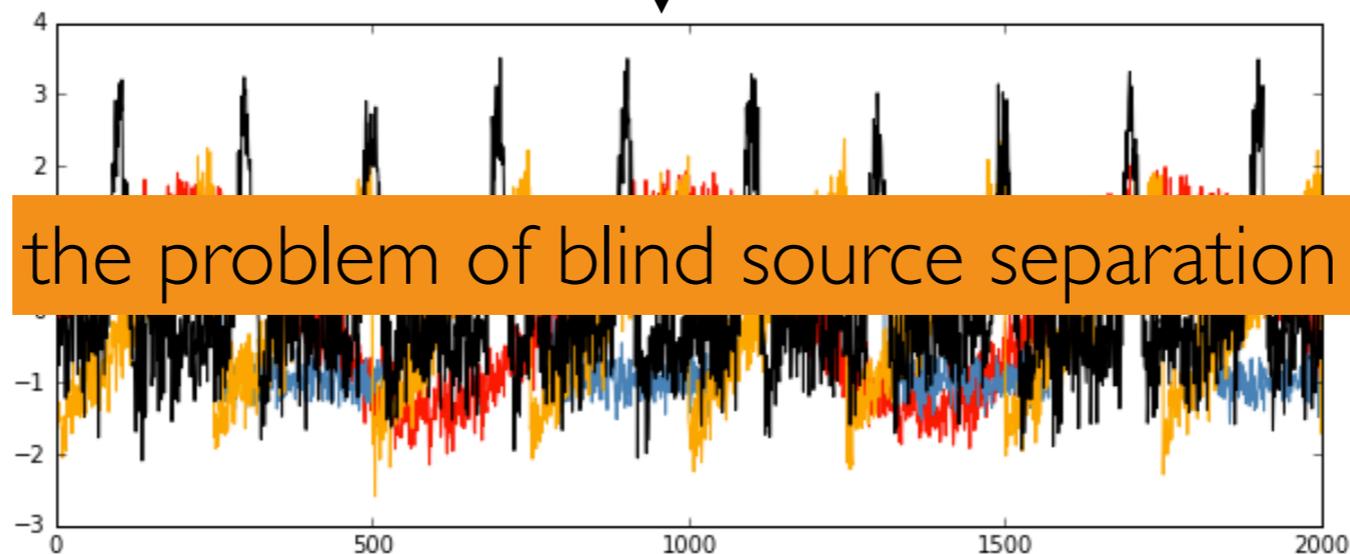
Signals



$$X = S \cdot A^T$$

How to recover the sources??

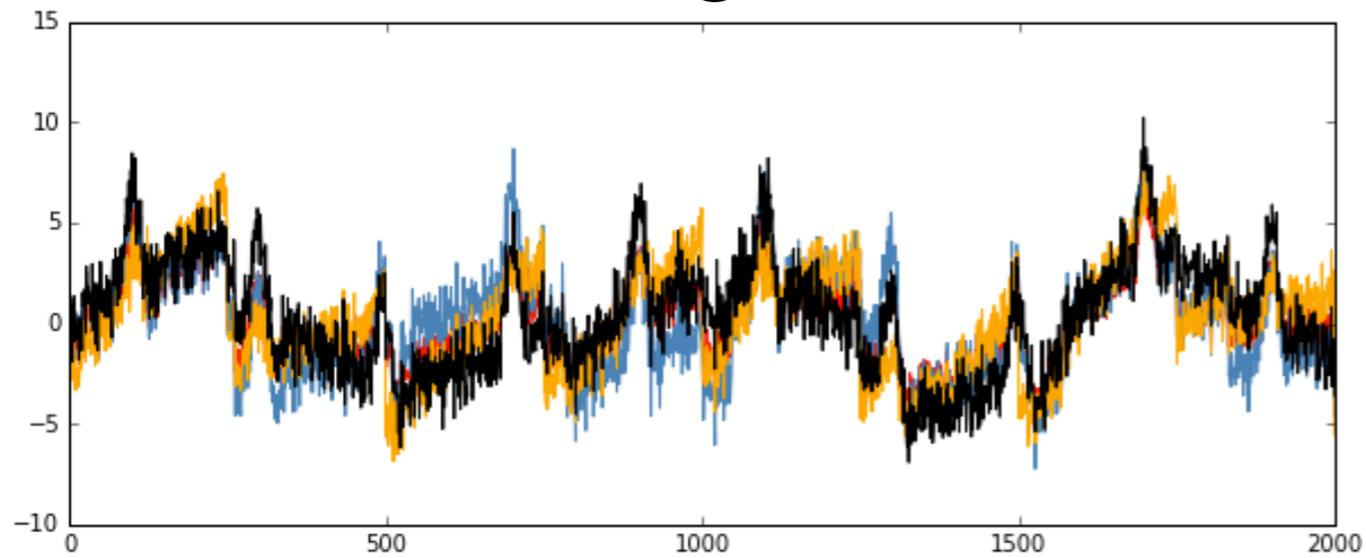
Problem: **S**, **A** unknown



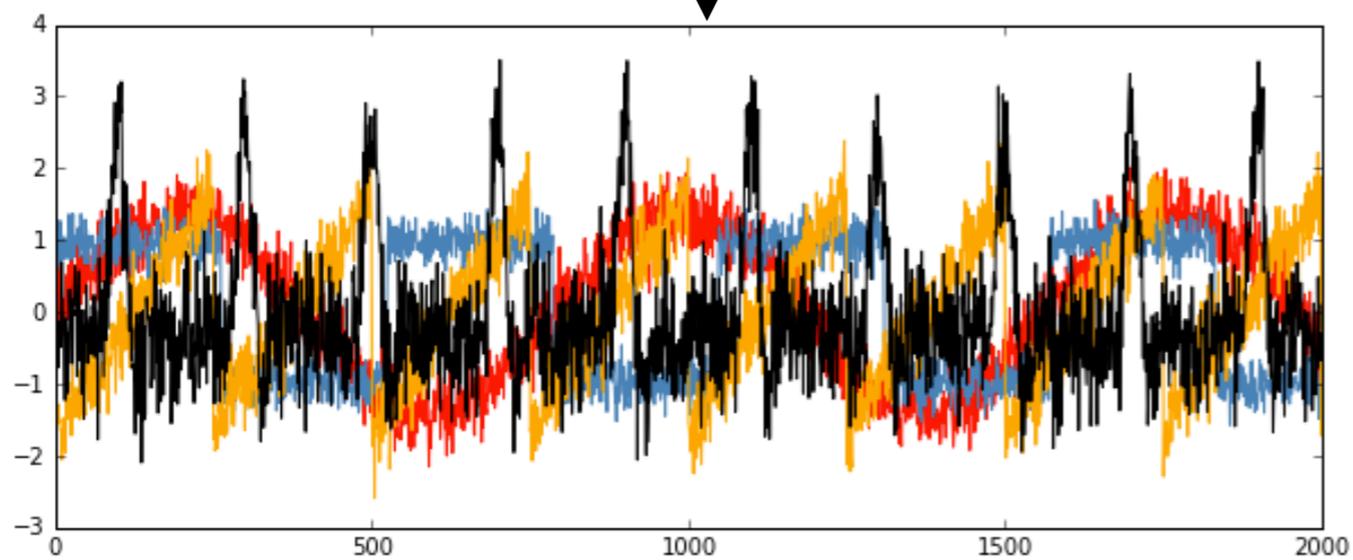
Sources

Independent component analysis

Signals



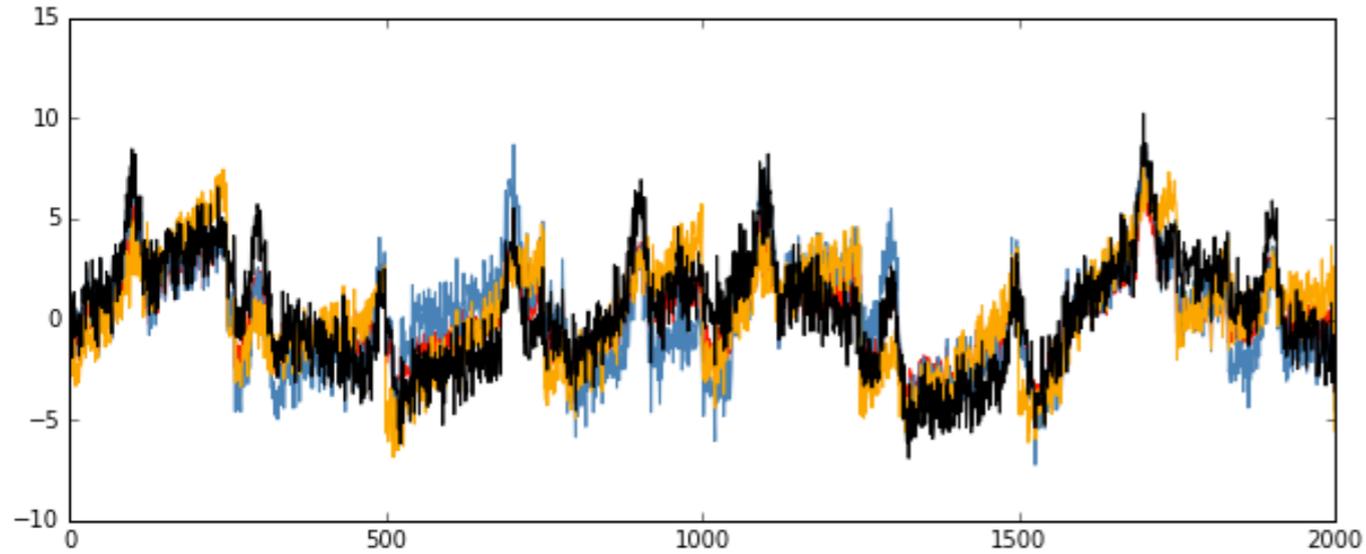
$$X = S \cdot A^T$$



Sources

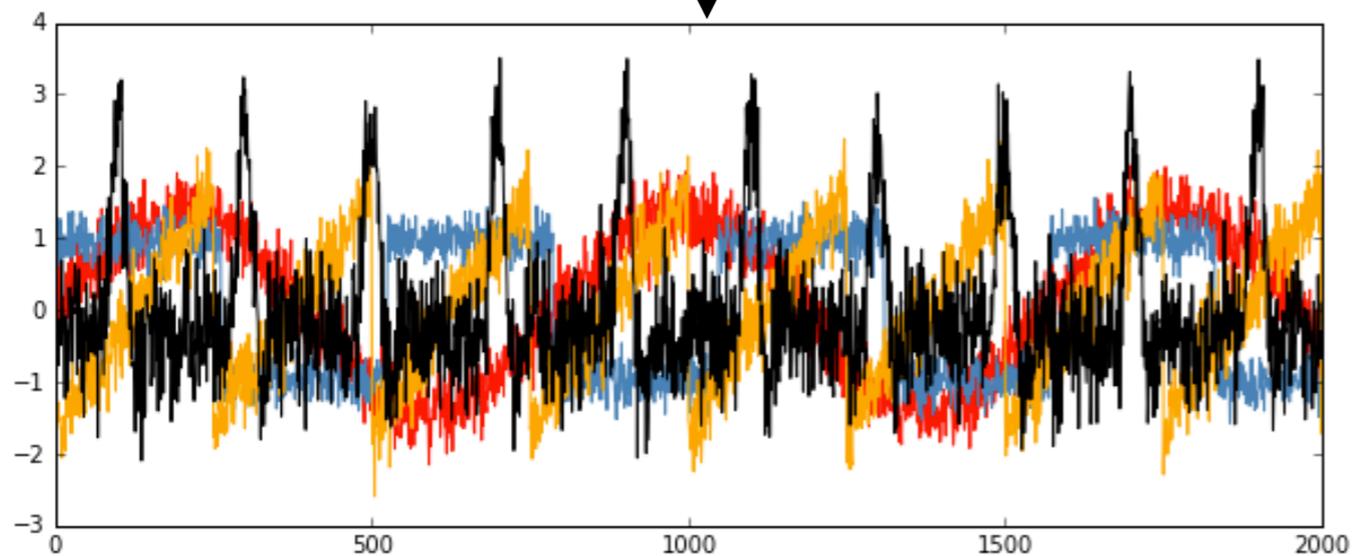
Independent component analysis

Signals



$$X = S \cdot A^T$$

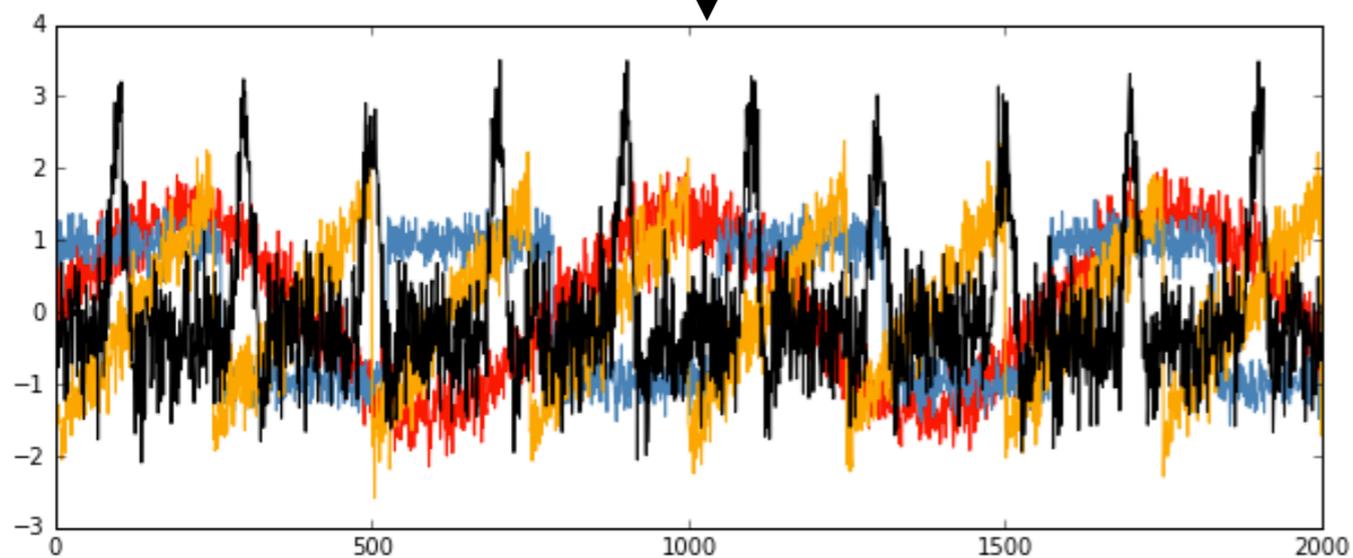
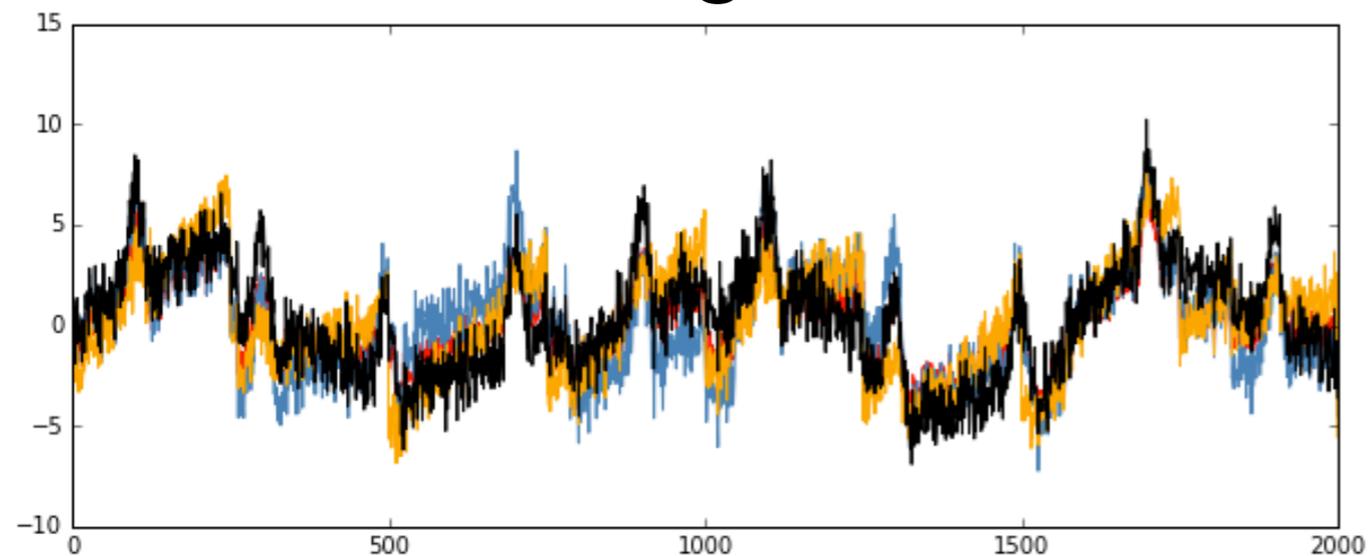
Core Idea: Mixtures are Gaussian ...



Sources

Independent component analysis

Signals



Sources

$$X = S \cdot A^T$$

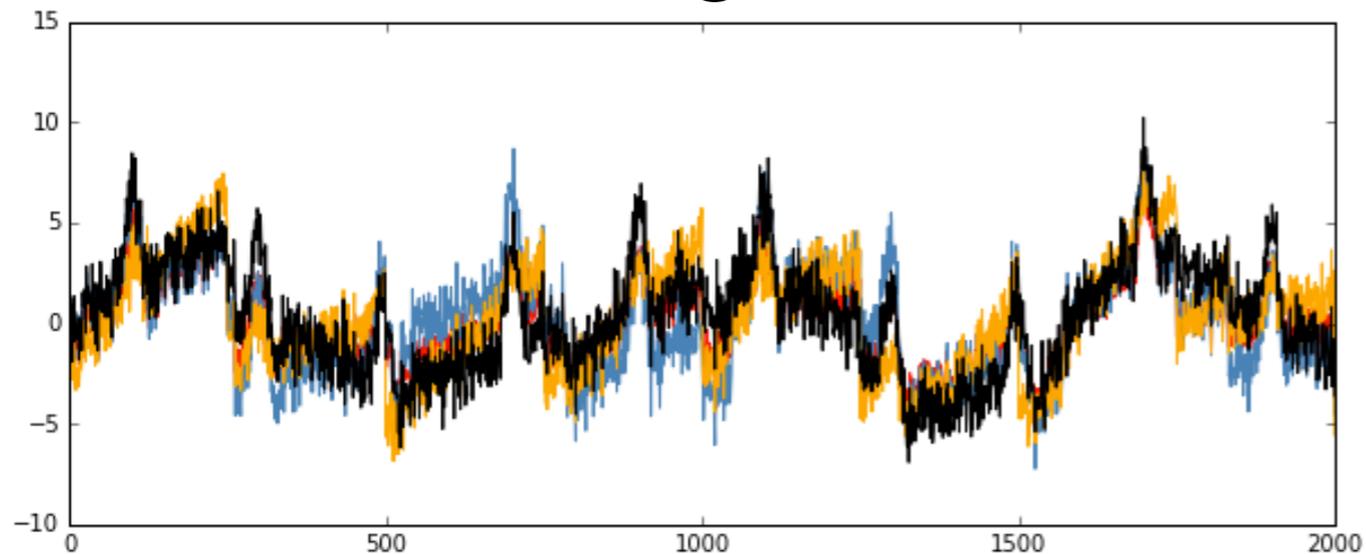
Core Idea: Mixtures are Gaussian ...

So sources should **not** be.

“Non-gaussian is independent.”

Independent component analysis

Signals

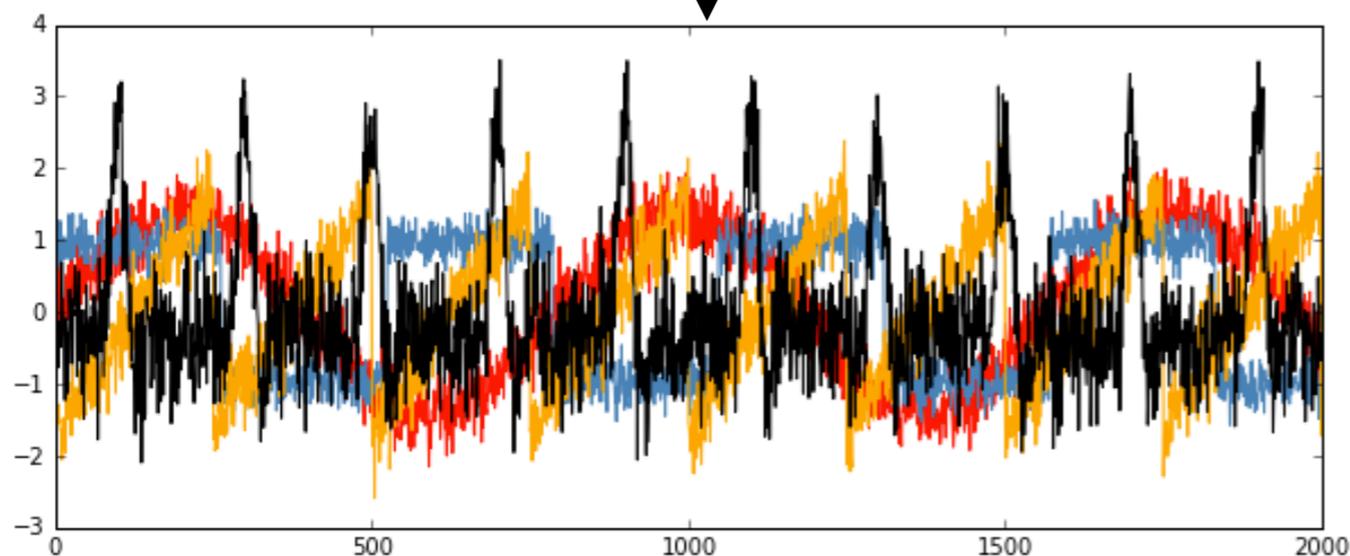


$$X = S \cdot A^T$$

Core Idea: Mixtures are Gaussian ...

Specifically: find \mathbf{w} that maximizes non-Gaussianity of

$$w^T X$$



So sources should **not** be.

“Non-gaussian is independent.”

Sources

Independent component analysis

Specifically: find \mathbf{w} that maximizes non-Gaussianity of $w^T X$

How to measure non-Gaussianity?

Independent component analysis

Specifically: find \mathbf{w} that maximizes non-Gaussianity of $w^T X$

How to measure non-Gaussianity?

a) Kurtosis of distribution

$$\text{kurt}(y) = \mathbb{E}[y^4] - 3\mathbb{E}(y^2)^2$$

Independent component analysis

Specifically: find \mathbf{w} that maximizes non-Gaussianity of $w^T X$

How to measure non-Gaussianity?

a) Kurtosis of distribution

$$\text{kurt}(y) = \mathbb{E}[y^4] - 3\mathbb{E}(y^2)^2$$

y is Gaussian \longrightarrow $\text{kurt}(y) = 0$

Independent component analysis

Specifically: find \mathbf{w} that maximizes non-Gaussianity of $w^T X$

How to measure non-Gaussianity?

a) Kurtosis of distribution

$$\text{kurt}(y) = \mathbb{E}[y^4] - 3\mathbb{E}(y^2)^2$$

$$y \text{ is Gaussian} \longrightarrow \text{kurt}(y) = 0$$

b) “Negentropy”

$$\text{entropy: } H(Y) = - \sum_i P(Y = a_i) \log P(Y = a_i)$$

Independent component analysis

Specifically: find \mathbf{w} that maximizes non-Gaussianity of $w^T X$

How to measure non-Gaussianity?

a) Kurtosis of distribution

$$\text{kurt}(y) = \mathbb{E}[y^4] - 3\mathbb{E}(y^2)^2$$

$$y \text{ is Gaussian} \longrightarrow \text{kurt}(y) = 0$$

b) “Negentropy”

$$\text{entropy: } H(Y) = - \sum_i P(Y = a_i) \log P(Y = a_i)$$

Key fact: entropy is maximized by Gaussian, fixing variance.

Independent component analysis

Specifically: find \mathbf{w} that maximizes non-Gaussianity of $w^T X$

How to measure non-Gaussianity?

a) Kurtosis of distribution

$$\text{kurt}(y) = \mathbb{E}[y^4] - 3\mathbb{E}(y^2)^2$$

$$y \text{ is Gaussian} \longrightarrow \text{kurt}(y) = 0$$

b) “Negentropy”

$$\text{entropy: } H(Y) = - \sum_i P(Y = a_i) \log P(Y = a_i)$$

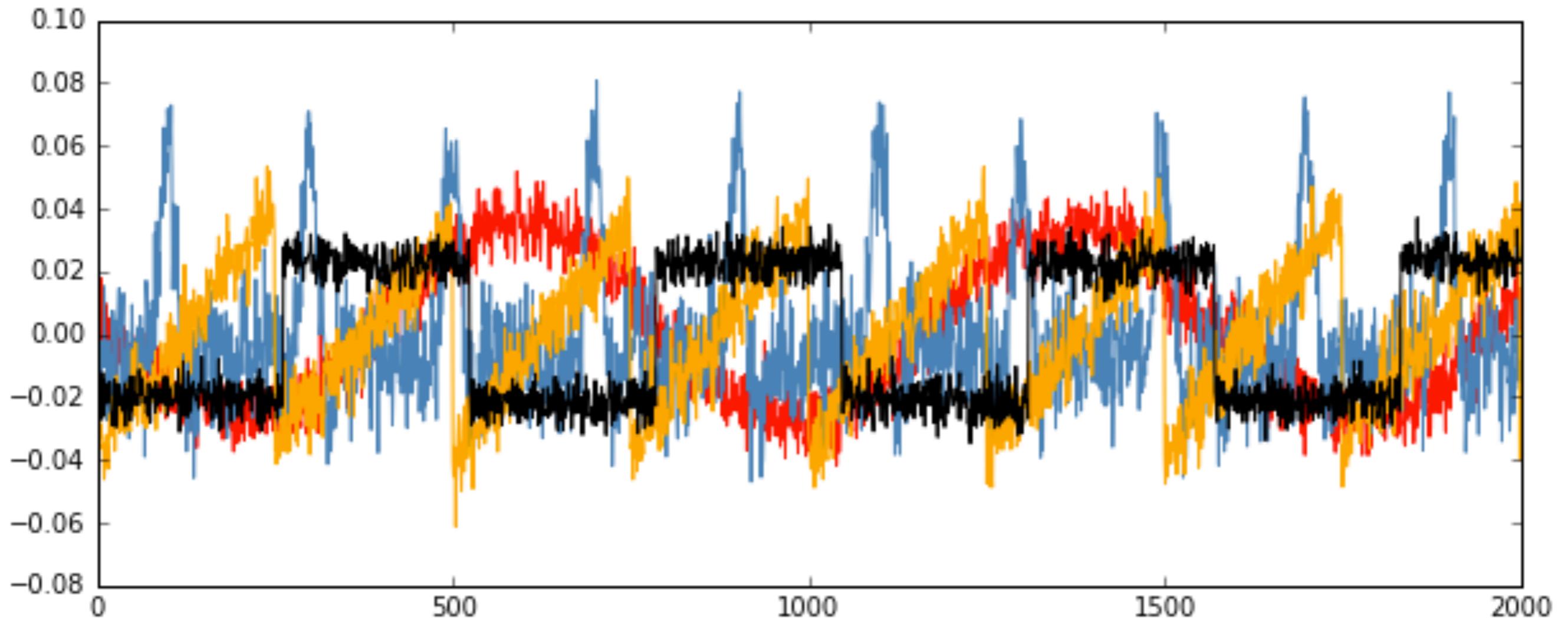
Key fact: entropy is maximized by Gaussian, fixing variance.

$$\text{negentropy: } J(y) = H(y_{\text{gaussian}}) - H(y)$$

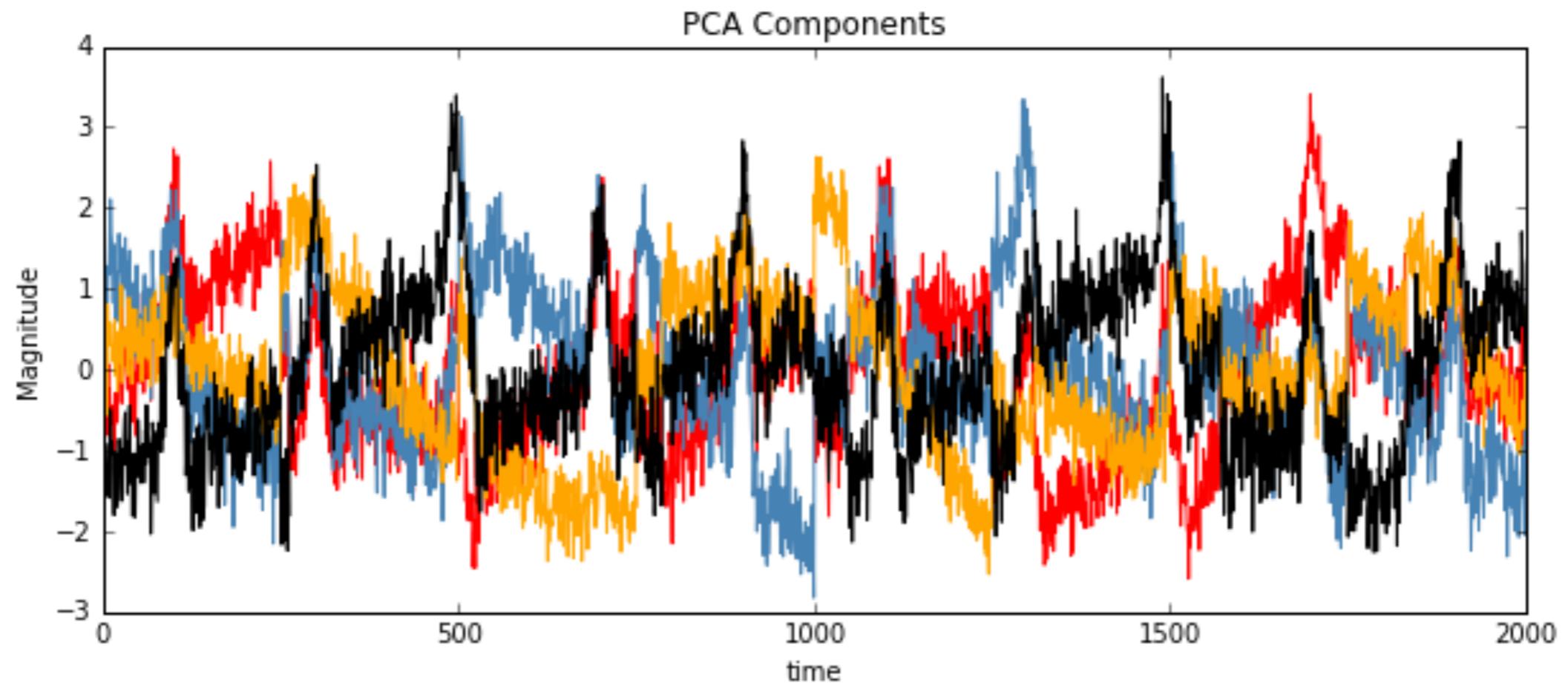
Independent component analysis

[Ipython: FastICA algorithm]

ICA results

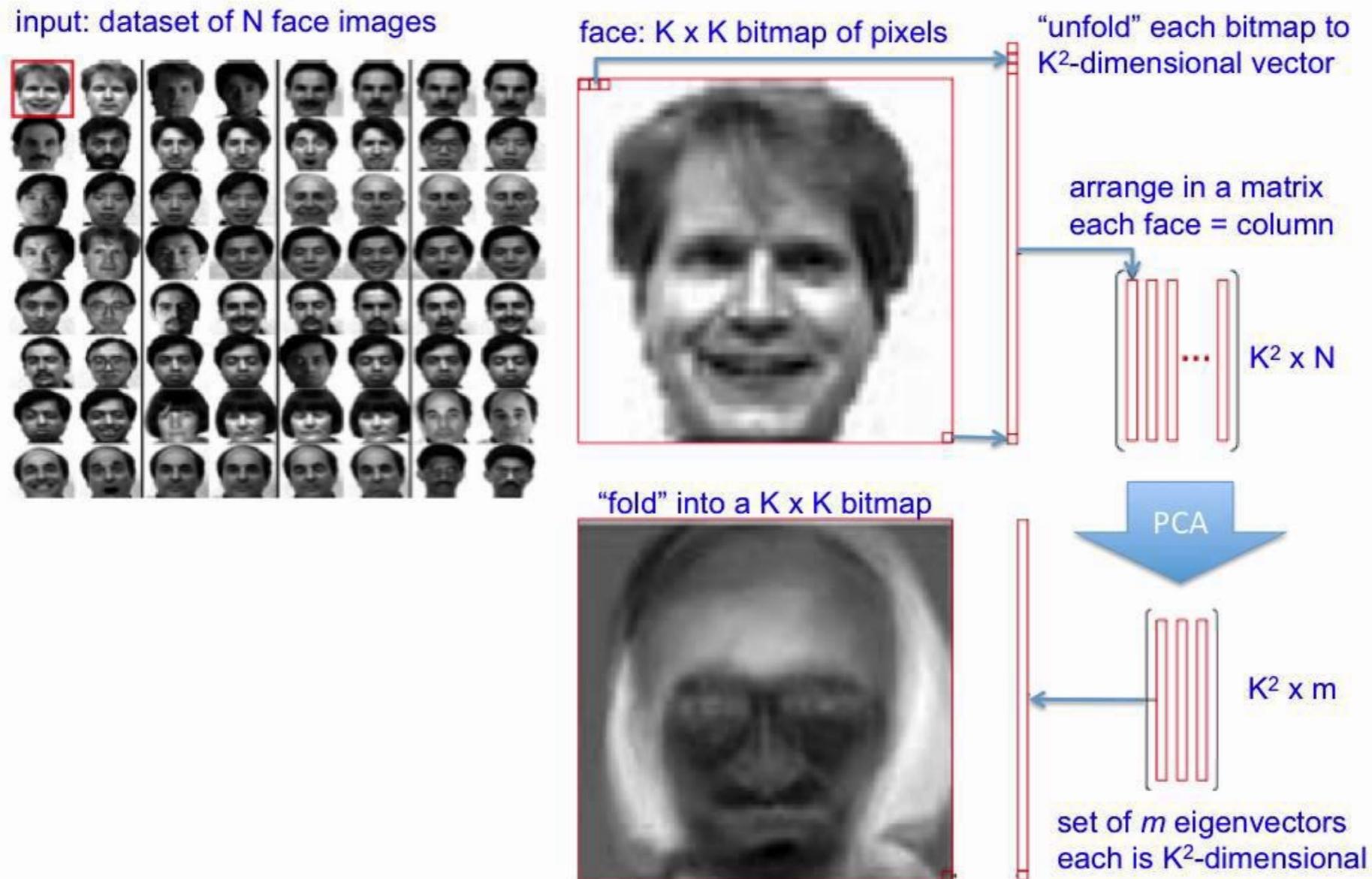


PCA results



PCA

PCA example: Eigen Faces

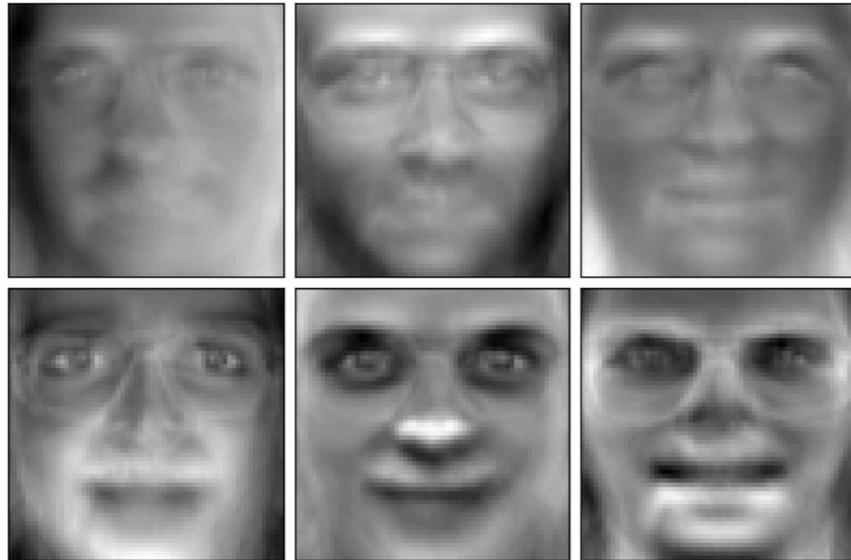


Independent component analysis

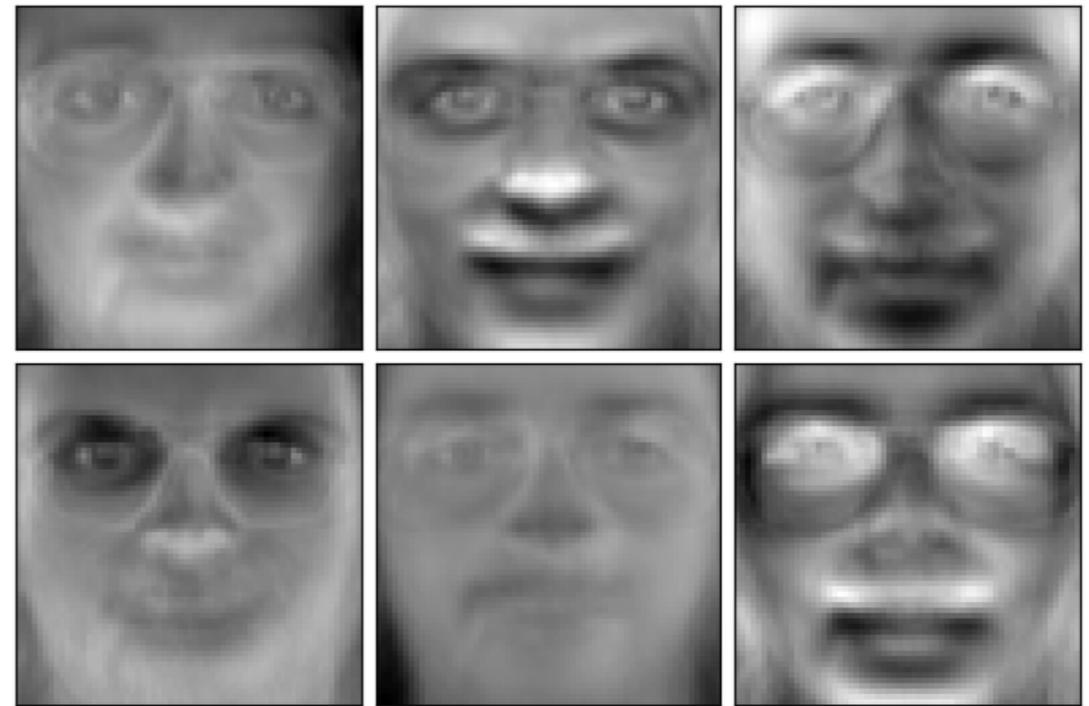
$$x = A \cdot s$$

where the s_i are statistically independent signals

genfaces - PCA using randomized SVD - Train time 0.1s



Independent components - FastICA - Train time 0.2s



First centered Olivetti faces



Independent component analysis

fMRI response data collected* on 165 commonly heard natural sound stimuli.

Man speaking
Flushing toilet
Pouring liquid
Tooth-brushing
Woman speaking
Car accelerating
Biting and chewing
Laughing
Typing
Car engine starting
Running water
Breathing
Keys jangling
Dishes clanking
Ringtone
Microwave
Dog barking

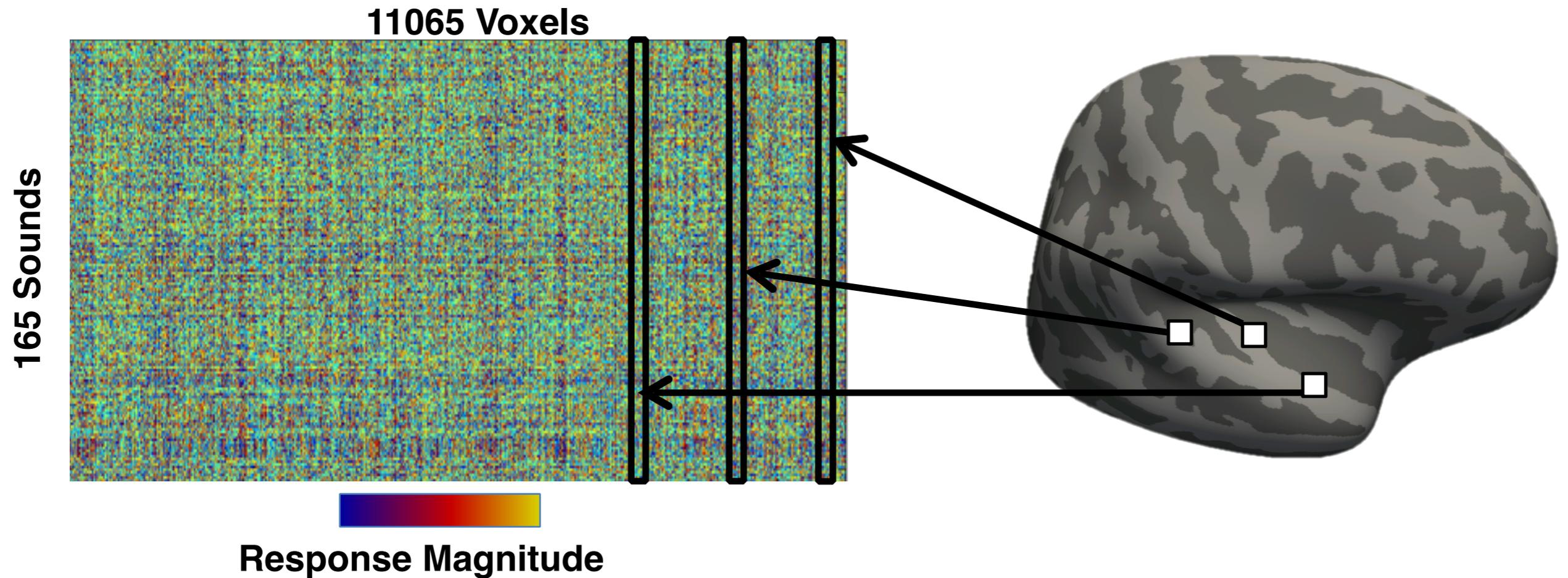
Road traffic
Zipper
Cellphone vibrating
Water dripping
Scratching
Car windows
Telephone ringing
Chopping food
Telephone dialing
Girl speaking
Car horn
Writing
Computer startup sound
Background speech
Songbird
Pouring water
Pop song
Water boiling

Guitar
Coughing
Crumpling paper
Siren
Splashing water
Computer speech
Alarm clock
Walking with heels
Vacuum
Wind
Boy speaking
Chair rolling
Rock song
Door knocking
●
●
●

*Sam Norman-Haignere, Nancy Kanwisher, and Josh McDermott

Independent component analysis

For each voxel, measured average response to each sound:



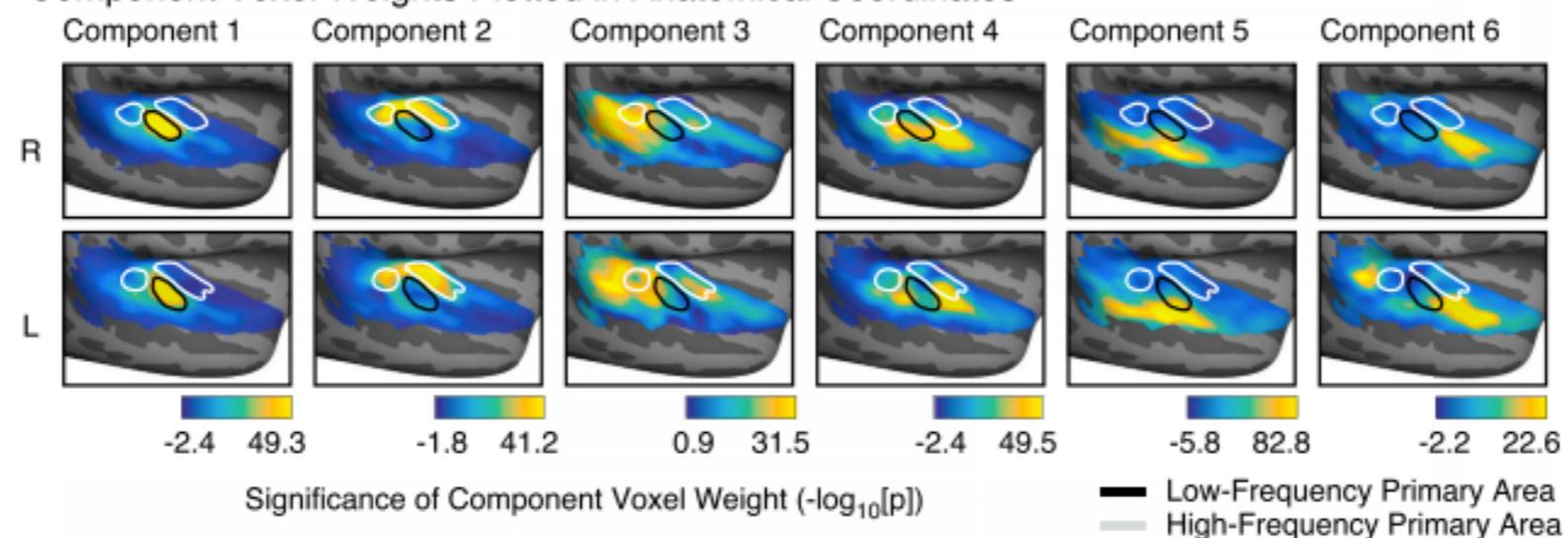
Data matrix: voxels \times sounds.

Independent component analysis

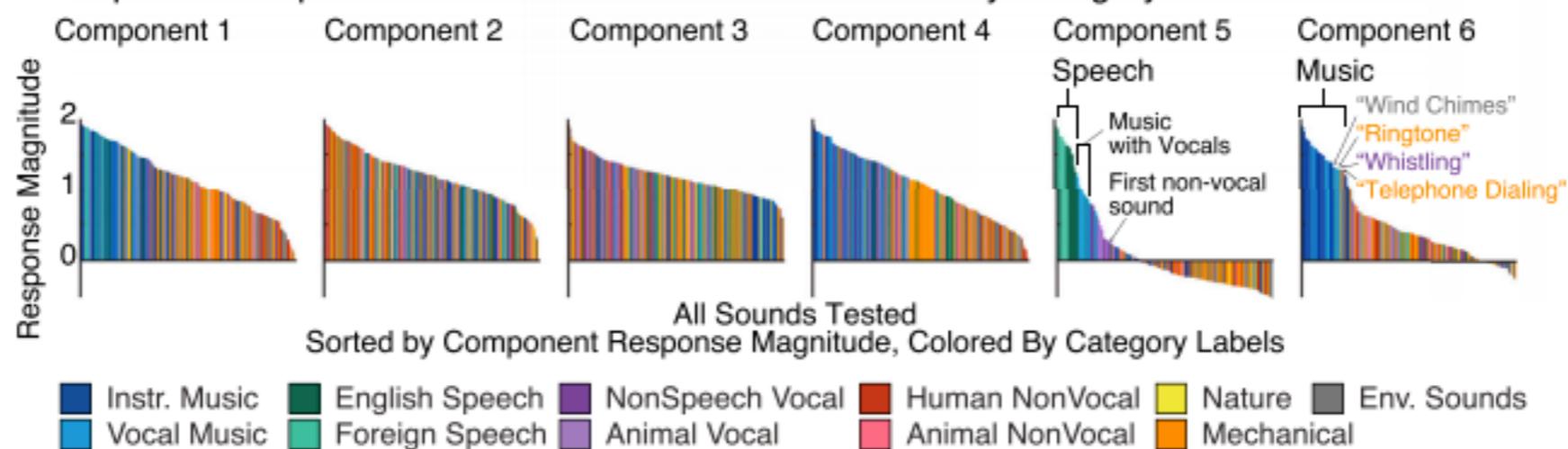
$$x = A \cdot s$$

where the s_i are statistically independent signals

B Component Voxel Weights Plotted in Anatomical Coordinates



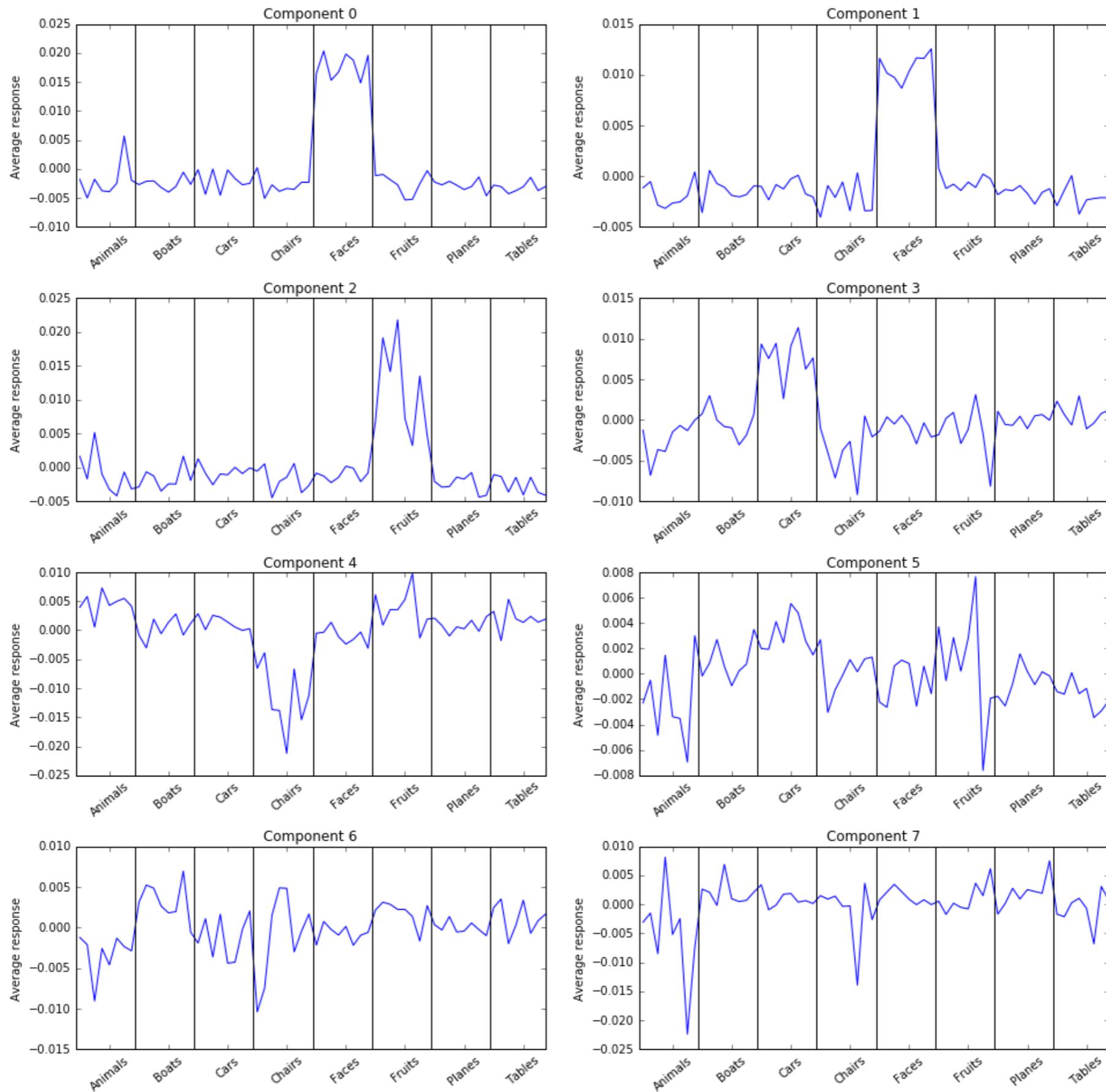
D Component Response Profiles to All 165 Sounds Colored by Category



Independent component analysis

[Ipython: ICA on neural data]

Independent component analysis



Non-negative Matrix Factorization (NMF)



Sebastian Seung

$$X \sim W \cdot H \quad W, H \geq 0$$

minimize: $\frac{1}{2} \sum_{ij} (X_{ij} - (W \cdot H)_{ij})^2$

Other Factorization Methods

Non-negative Matrix Factorization (NMF)

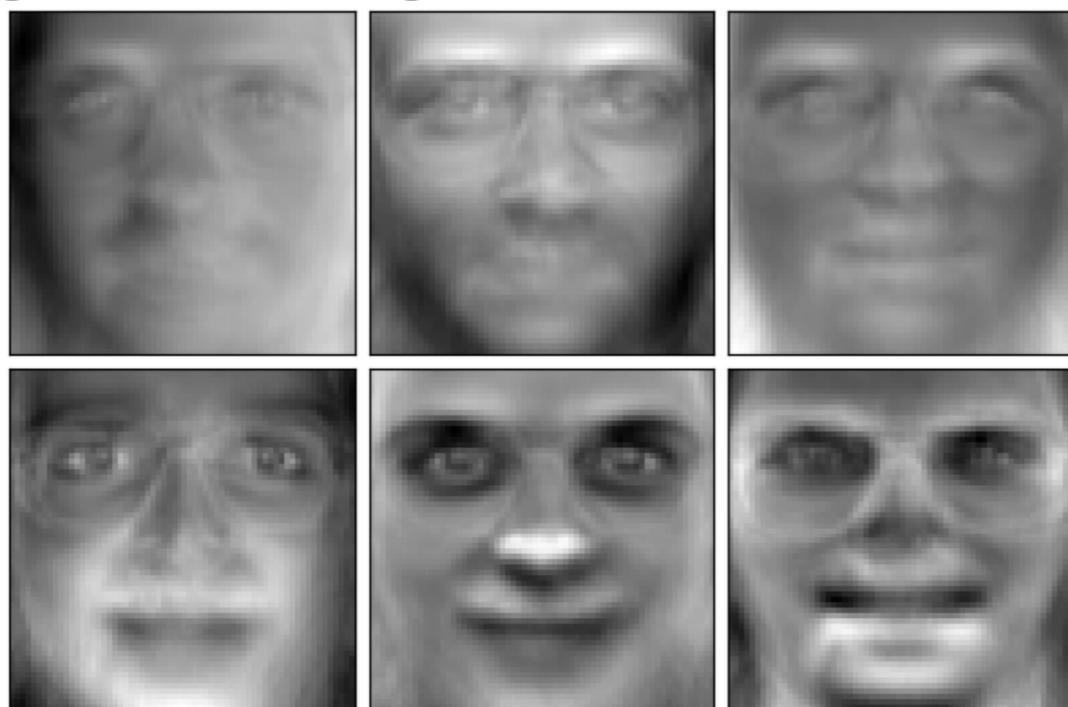


Sebastian Seung

$$X \sim W \cdot H \quad W, H \geq 0$$

minimize: $\frac{1}{2} \sum_{ij} (X_{ij} - (W \cdot H)_{ij})^2$

genfaces - PCA using randomized SVD - Train time 0.1s



Non-negative components - NMF - Train time 0.3s



Other Factorization Methods

Non-negative Matrix Factorization (NMF)



Sebastian Seung

$$X \sim W \cdot H \quad W, H \geq 0$$

minimize:
$$\frac{1}{2} \sum_{ij} (X_{ij} - (W \cdot H)_{ij})^2$$

genfaces - PCA using randomized SVD - Train time 0.1s



Non-negative components - NMF - Train time 0.3s



regularization:

$$\frac{1}{2} \sum_{ij} (X_{ij} - (W \cdot H)_{ij})^2 + \alpha(\|W\|_1 + \|H\|_1)$$

Other Factorization Methods

Non-negative Matrix Factorization (NMF)

NEURORESOURCE

Simultaneous Denoising, Deconvolution, and Demixing of Calcium Imaging Data

Eftychios A. Pnevmatikakis[✉], Daniel Soudry, Yuanjun Gao, Timothy A. Machado, Josh Merel, David Pfau, Thomas Reardon, Yu Mu, Clay Lacefield, Weijian Yang, Misha Ahrens, Randy Bruno, Thomas M. Jessell, Darcy S. Peterka, Rafael Yuste, Liam Paninski[✉]

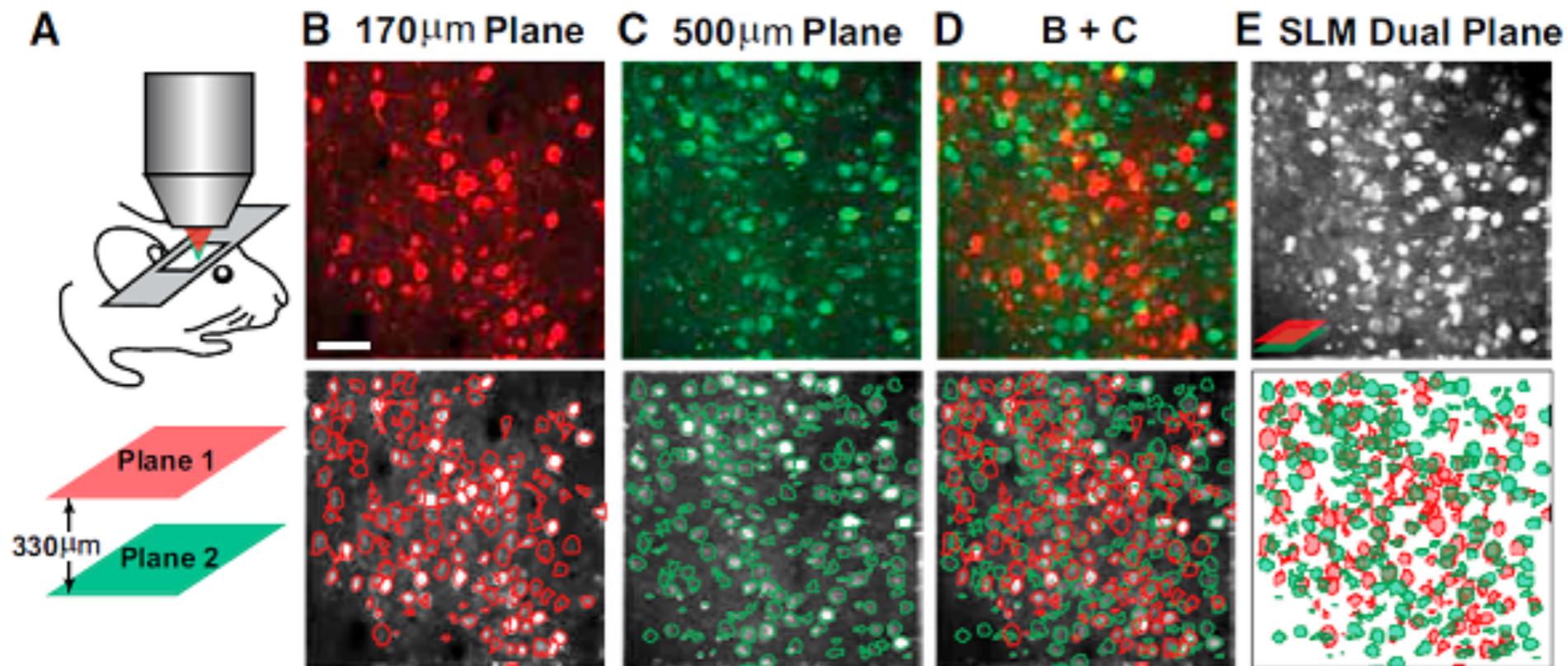
Published Online: January 07, 2016

Open Archive



DOI: <http://dx.doi.org/10.1016/j.neuron.2015.11.037> | CrossMark

Article Info

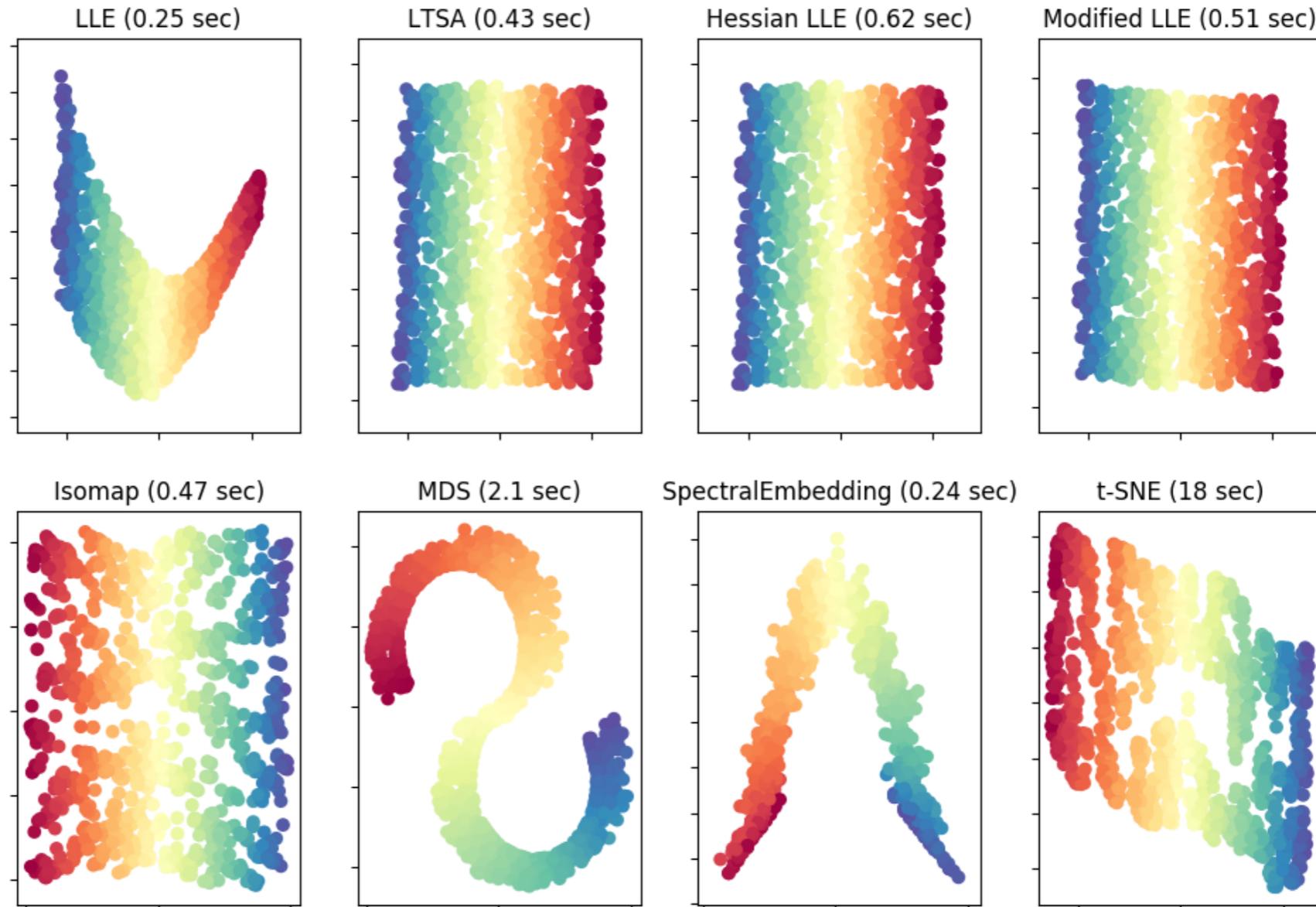
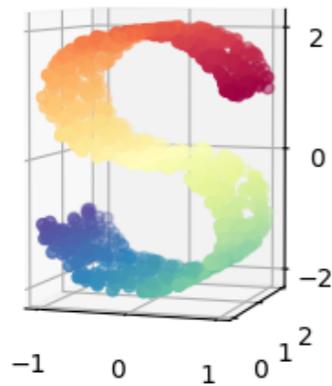


Manifold Learning

Arrange the data so that values fill out “nice” coordinates on a “nicely-shaped” manifold

Manifold Learning with 1000 points, 10 neighbors

True 3D data



Various embeddings in 2D

Multi-Dimensional Scaling (MDS)

Goal: find an “embedding function”

$$F : \mathbb{R}^k \longrightarrow \mathbb{R}^l$$

Multi-Dimensional Scaling (MDS)

Goal: find an “embedding function”

$$F : \mathbb{R}^k \longrightarrow \mathbb{R}^l$$

where dimension has been dramatically reduced:

$$l \ll k$$

Multi-Dimensional Scaling (MDS)

Goal: find an “embedding function”

$$F : \mathbb{R}^k \longrightarrow \mathbb{R}^l$$

where dimension has been dramatically reduced:

$$l \ll k$$

but such that distances have been preserved:

$$d(F(x_i), F(x_j)) \sim d(x_i, x_j)$$

Manifold Learning

Multi-Dimensional Scaling (MDS)

Goal: find an “embedding function”

$$F : \mathbb{R}^k \longrightarrow \mathbb{R}^l$$

where dimension has been dramatically reduced:

$$l \ll k$$

but such that distances have been preserved:

$$d(F(x_i), F(x_j)) \sim d(x_i, x_j)$$

distance in reduced-dim embedding

distance in original space

Manifold Learning

Multi-Dimensional Scaling (MDS)

$$F : \mathbb{R}^k \longrightarrow \mathbb{R}^l; l \ll k; d(F(x_i), F(x_j)) \sim d(x_i, x_j)$$

Loss function:

$$\text{Stress}(F(x_1), \dots, F(x_N)) = \left(\sum_{i \neq j} (d(x_i, x_j) - \|F(x_i) - F(x_j)\|)^2 \right)^{\frac{1}{2}}$$

Optimization by seeking to minimize stress.

Manifold Learning

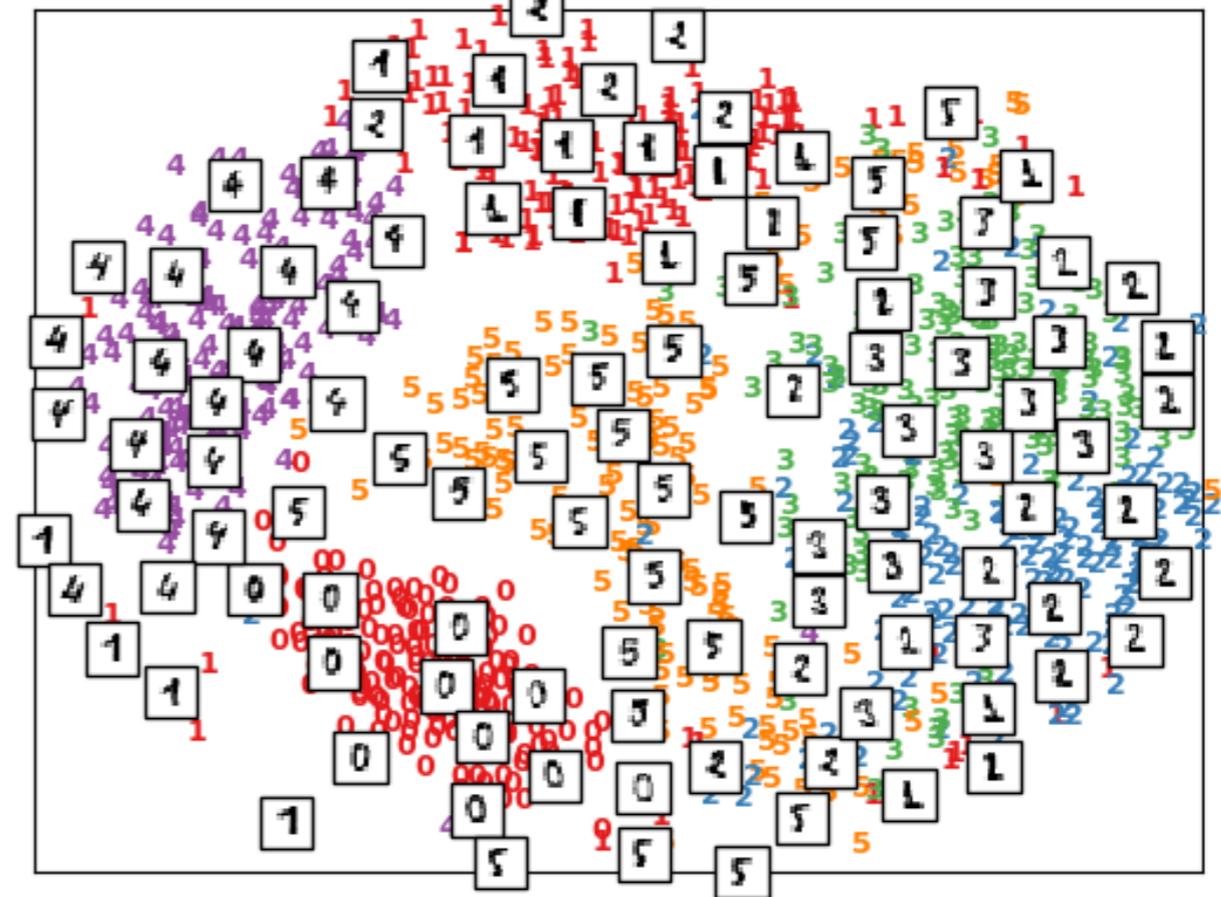
Multi-Dimensional Scaling (MDS)

$$F : \mathbb{R}^k \longrightarrow \mathbb{R}^l; l \ll k; d(F(x_i), F(x_j)) \sim d(x_i, x_j)$$

A selection from the 64-dimensional digits dataset

0	1	2	3	4	5	0	1	2	3	4	5	0	1	2	3	4	5	0	5
5	5	0	4	1	3	5	1	0	0	2	2	2	0	1	2	3	3	3	3
4	4	1	5	0	5	2	2	0	0	1	3	2	1	4	3	1	3	1	4
3	4	4	0	5	3	1	5	4	4	2	2	2	5	5	4	4	0	0	1
2	3	4	5	0	1	2	3	4	5	0	1	2	3	4	5	0	5	5	5
0	4	4	3	5	1	0	0	2	2	2	0	4	2	3	3	3	4	4	4
1	5	0	5	2	2	0	0	1	3	2	1	3	1	3	4	4	3	1	4
0	5	3	4	5	4	4	1	2	2	5	5	4	4	0	0	1	2	3	4
5	0	1	2	3	4	5	0	4	2	3	4	5	0	5	5	5	0	4	1
3	5	1	0	0	2	2	2	0	1	2	3	3	3	3	4	4	1	5	0
5	2	2	0	0	1	3	2	1	4	3	1	4	3	1	4	3	1	4	0
3	1	5	4	4	2	2	2	5	5	4	4	0	3	0	1	2	3	4	5
0	1	2	3	4	5	0	1	2	3	4	5	0	5	5	5	0	4	1	3
5	1	0	0	1	2	2	0	1	2	3	3	3	3	4	4	1	5	0	5
1	2	0	0	1	3	2	1	4	3	1	3	1	4	3	1	4	0	5	3
1	5	4	4	2	2	2	5	5	4	4	0	0	1	2	3	4	5	0	1
2	3	4	5	0	1	2	3	4	5	0	5	5	5	0	4	1	3	5	4
0	0	2	2	2	0	1	2	3	3	3	3	4	4	1	5	0	5	2	2
0	0	1	3	2	1	4	3	1	3	1	4	3	1	4	0	5	3	1	5
4	4	2	2	1	5	5	4	4	0	0	1	2	3	4	5	0	1	2	3

MDS embedding of the digits (time 2.72s)

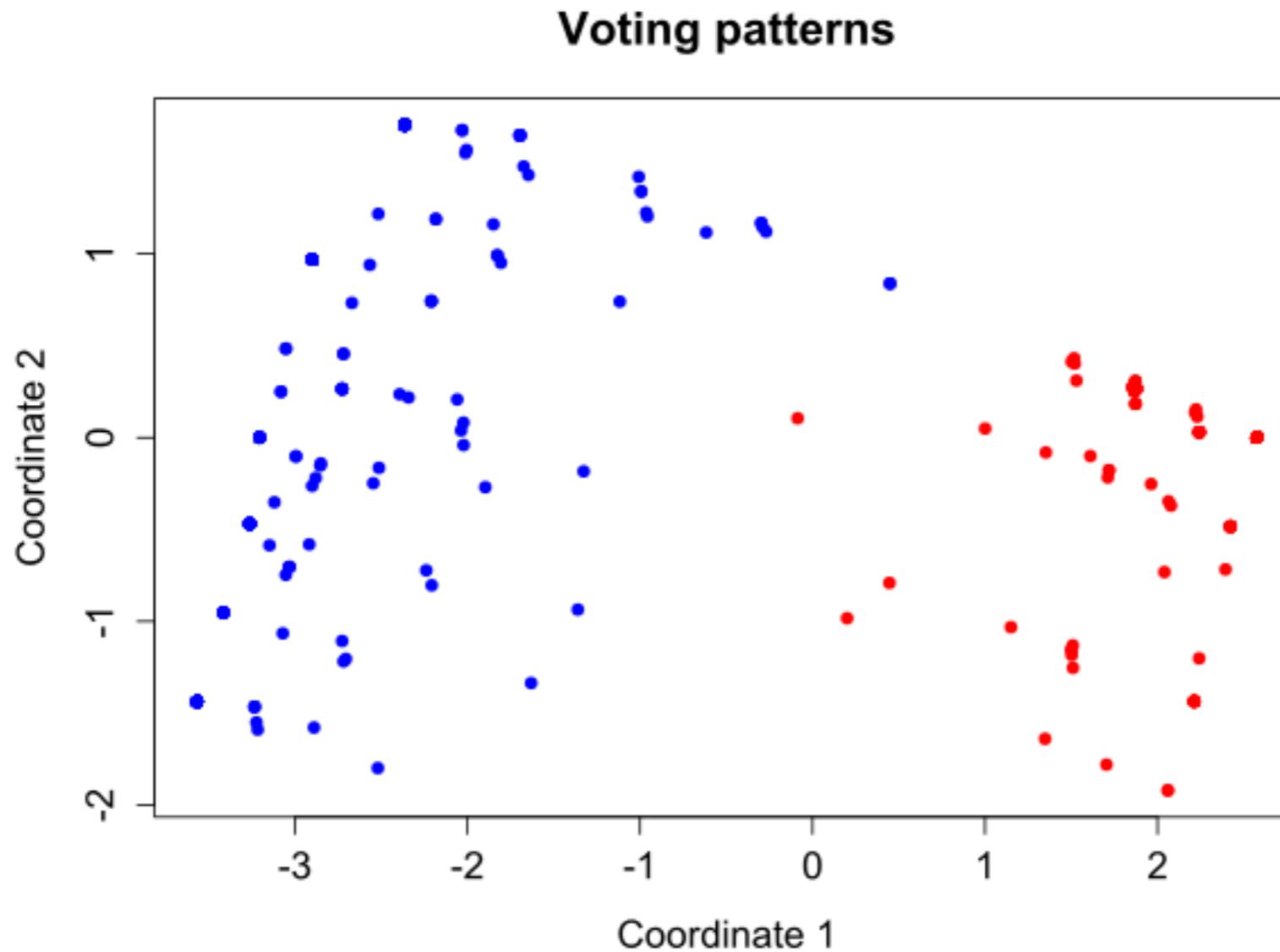


number of dimensions / learned via supervision.

Manifold Learning

Multi-Dimensional Scaling (MDS)

$$F : \mathbb{R}^k \longrightarrow \mathbb{R}^l; l \ll k; d(F(x_i), F(x_j)) \sim d(x_i, x_j)$$



“An example of classical multidimensional scaling applied to voting patterns in the United States House of Representatives. Each red dot represents one Republican member of the House, and each blue dot one Democrat.”

from Wikipedia

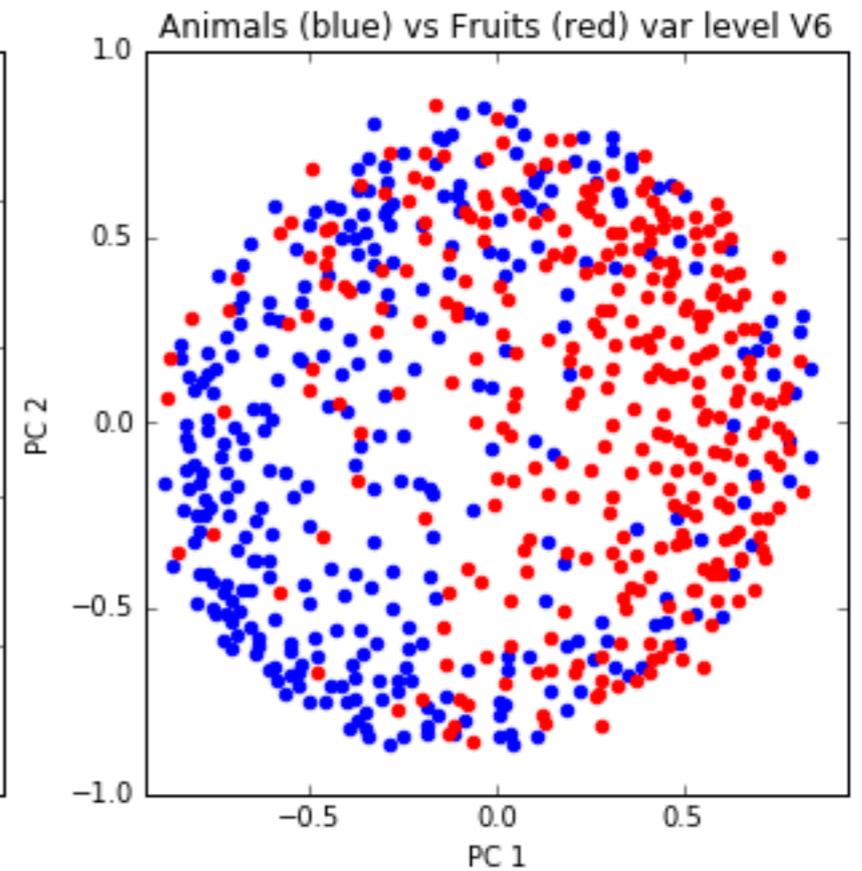
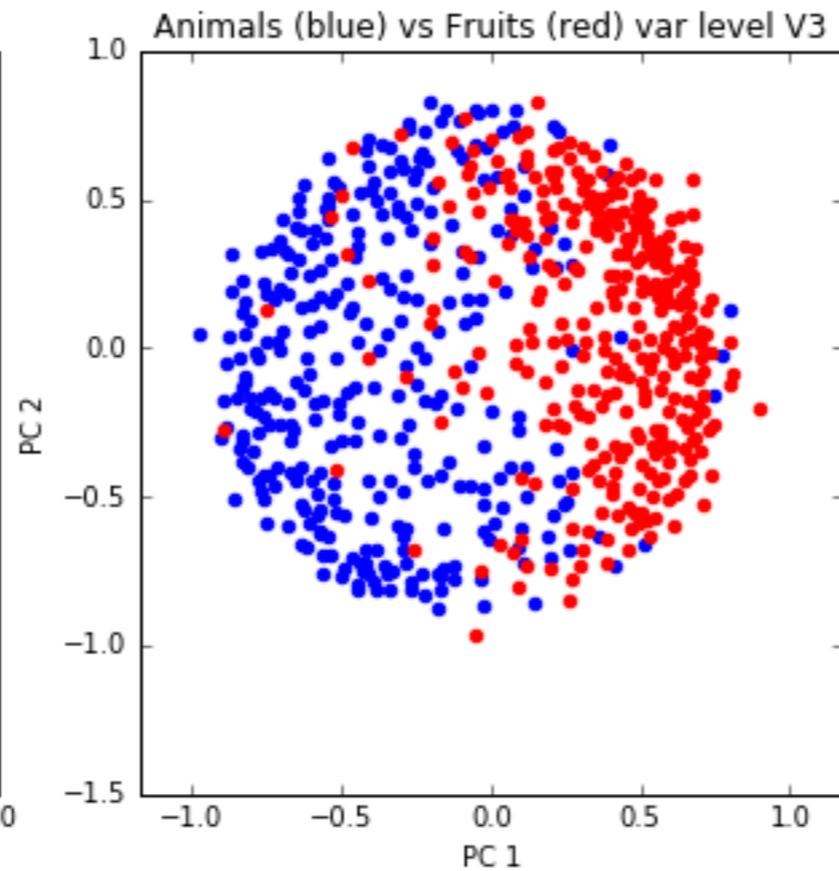
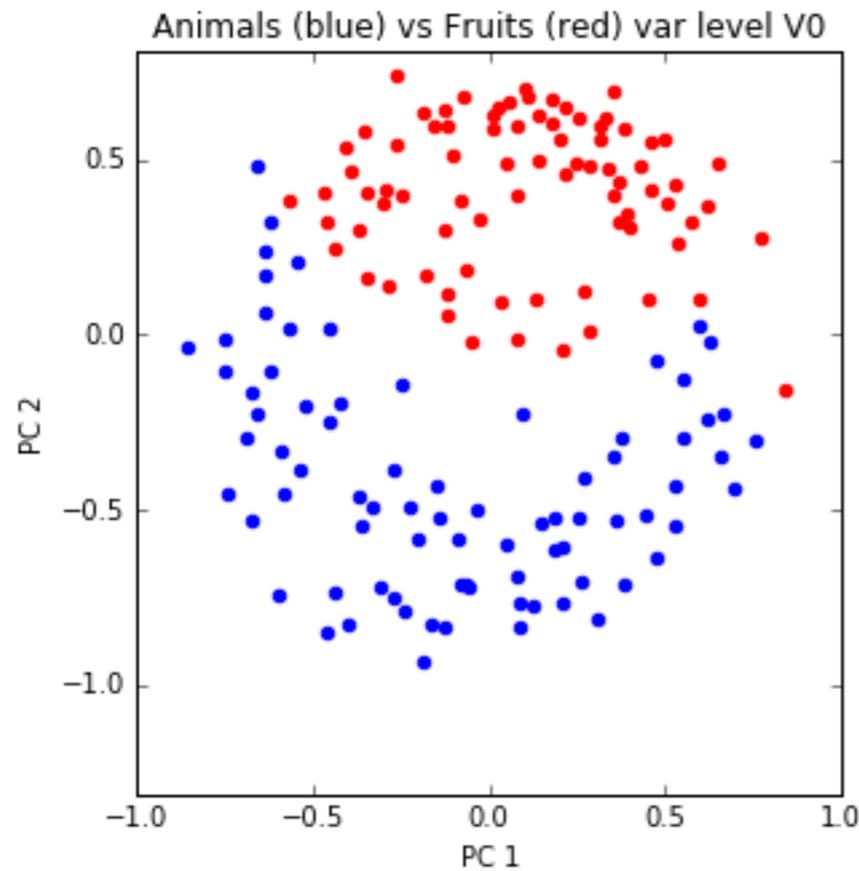
Independent component analysis

[Ipython: MDS on neural data]

Manifold Learning

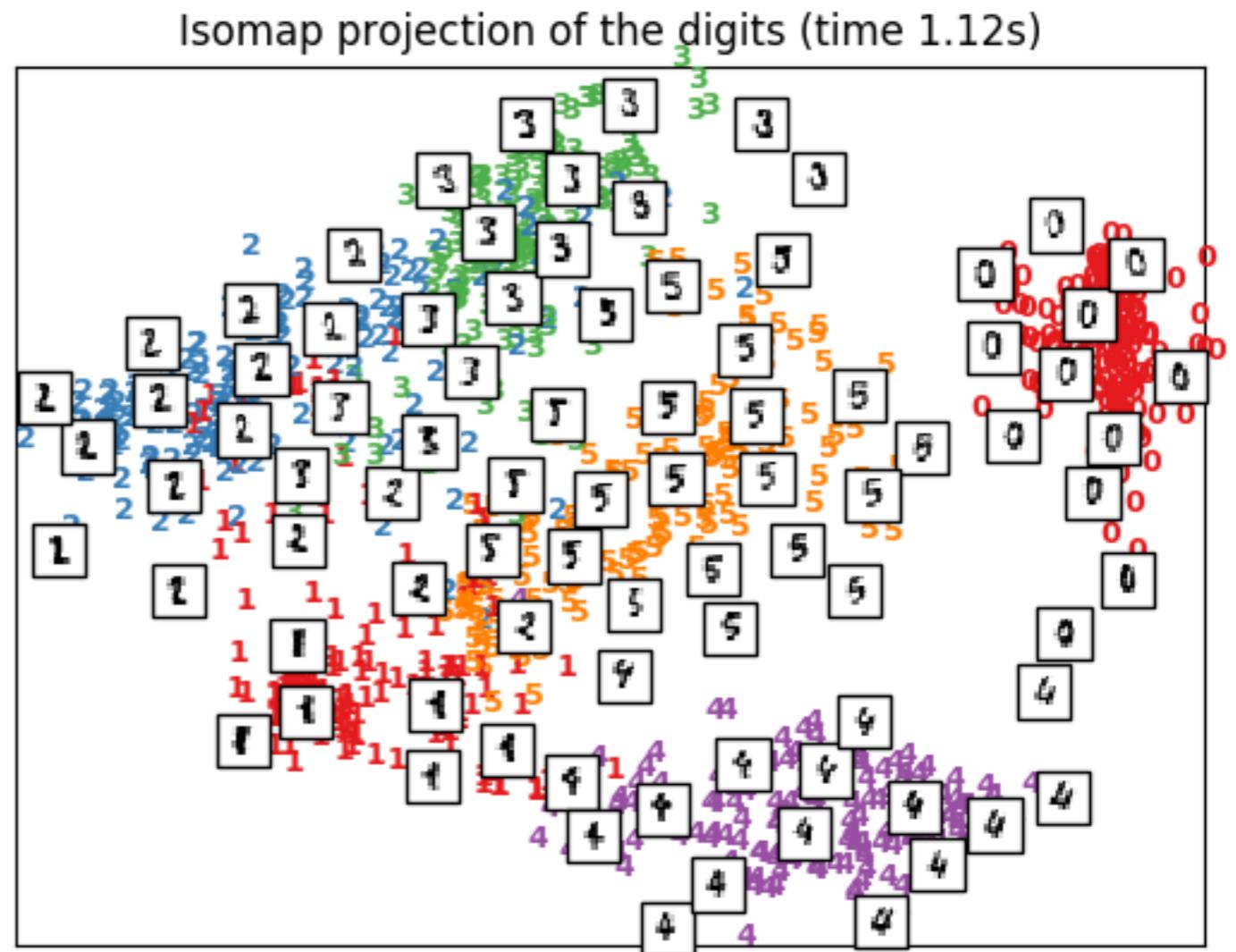
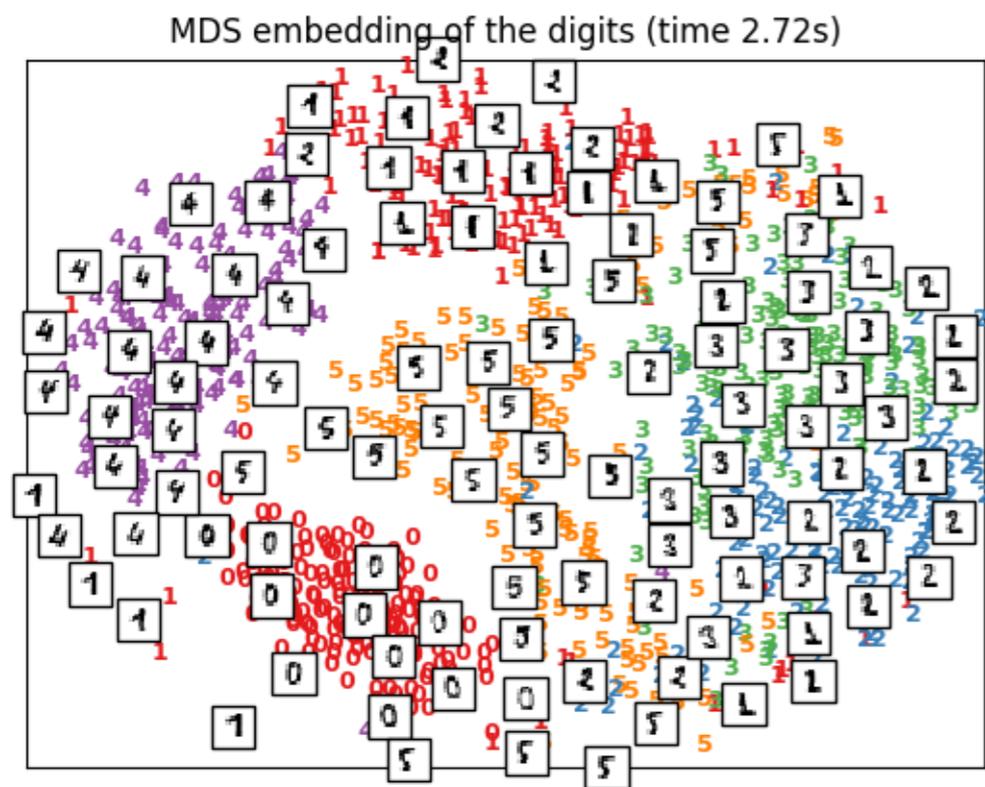
Multi-Dimensional Scaling (MDS)

$$F : \mathbb{R}^k \longrightarrow \mathbb{R}^l; l \ll k; d(F(x_i), F(x_j)) \sim d(x_i, x_j)$$



Manifold Learning

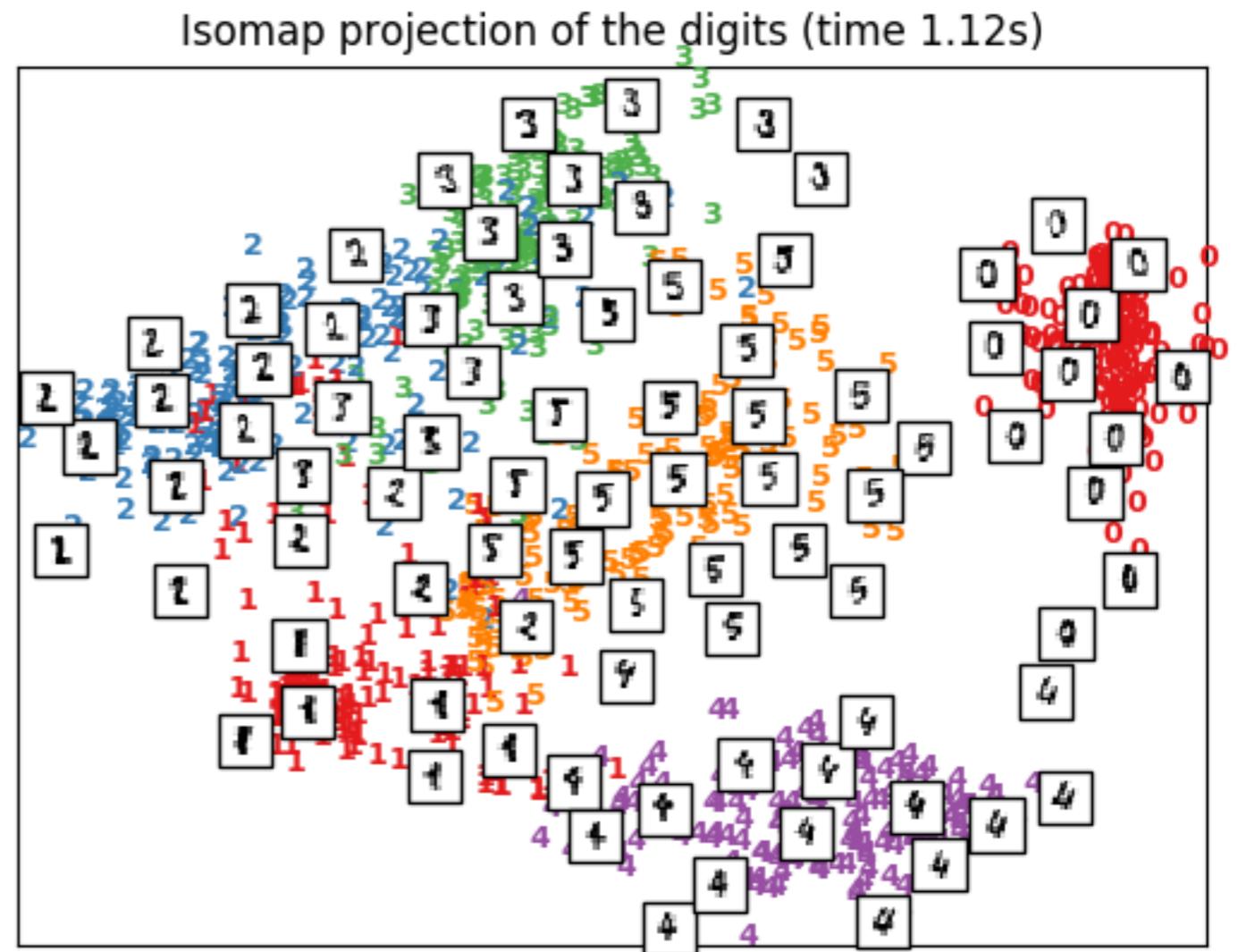
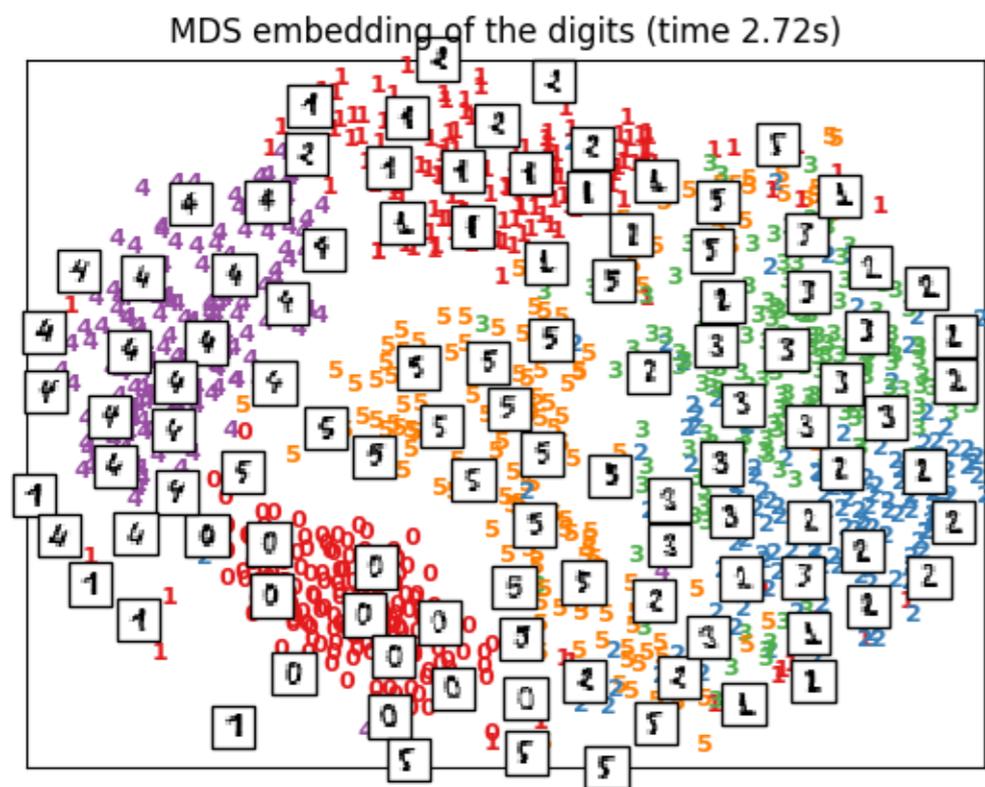
IsoMap: extends MDS by using computing distances relative to k -nearest neighbor graph in the high-dimensional space (e.g “manifold geodesics” rather than straight-line distance).



number of dimensions l and neighborhood k learned via supervision.

Manifold Learning

IsoMap: extends MDS by using computing distances relative to k -nearest neighbor graph in the high-dimensional space (e.g “manifold geodesics” rather than straight-line distance).



number of dimensions l and neighborhood k learned via supervision.

t-SNE chooses a particular generative model for data and for the compression

$$x \in \mathbb{R}^k; y \in \mathbb{R}^l; l \ll k; F(x) \mapsto y$$

t-SNE chooses a particular generative model for data and for the compression

$$x \in \mathbb{R}^k; y \in \mathbb{R}^l; l \ll k; F(x) \mapsto y$$

$$d(x_i, x_j) = \frac{\exp(-|x_i - x_j|^2 / 2\sigma_i^2)}{\sum_{k \neq i} \exp(-|x_i - x_k|^2 / 2\sigma_i^2)} + \text{symmetric term}$$

as if generated from gaussian joint probability

t-SNE chooses a particular generative model for data and for the compression

$$x \in \mathbb{R}^k; y \in \mathbb{R}^l; l \ll k; F(x) \mapsto y$$

$$d(x_i, x_j) = \frac{\exp(-|x_i - x_j|^2 / 2\sigma_i^2)}{\sum_{k \neq i} \exp(-|x_i - x_k|^2 / 2\sigma_i^2)} + \text{symmetric term}$$

as if generated from gaussian joint probability

$$d(y_i, y_j) = \frac{1 / (1 + (y_i - y_j)^2)}{\sum_{k \neq i} 1 / (1 + (y_i - y_k)^2)}$$

as if generated from student t distribution

Manifold Learning

t-SNE chooses a particular *generative model* for data and for the compression

$$x \in \mathbb{R}^k; y \in \mathbb{R}^l; l \ll k; F(x) \mapsto y$$

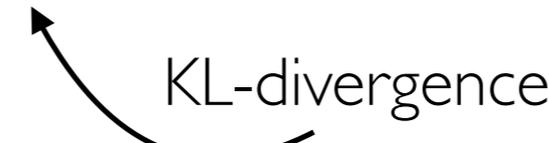
$$d(x_i, x_j) = \frac{\exp(-|x_i - x_j|^2 / 2\sigma_i^2)}{\sum_{k \neq i} \exp(-|x_i - x_k|^2 / 2\sigma_i^2)} + \text{symmetric term}$$

as if generated from gaussian joint probability

$$d(y_i, y_j) = \frac{1 / (1 + (y_i - y_j)^2)}{\sum_{k \neq i} 1 / (1 + (y_i - y_k)^2)}$$

as if generated from student t distribution

minimize $\sum_{i,j} d(x_i, x_j) \log \left[\frac{d(x_i, x_j)}{d(y_i, y_j)} \right]$



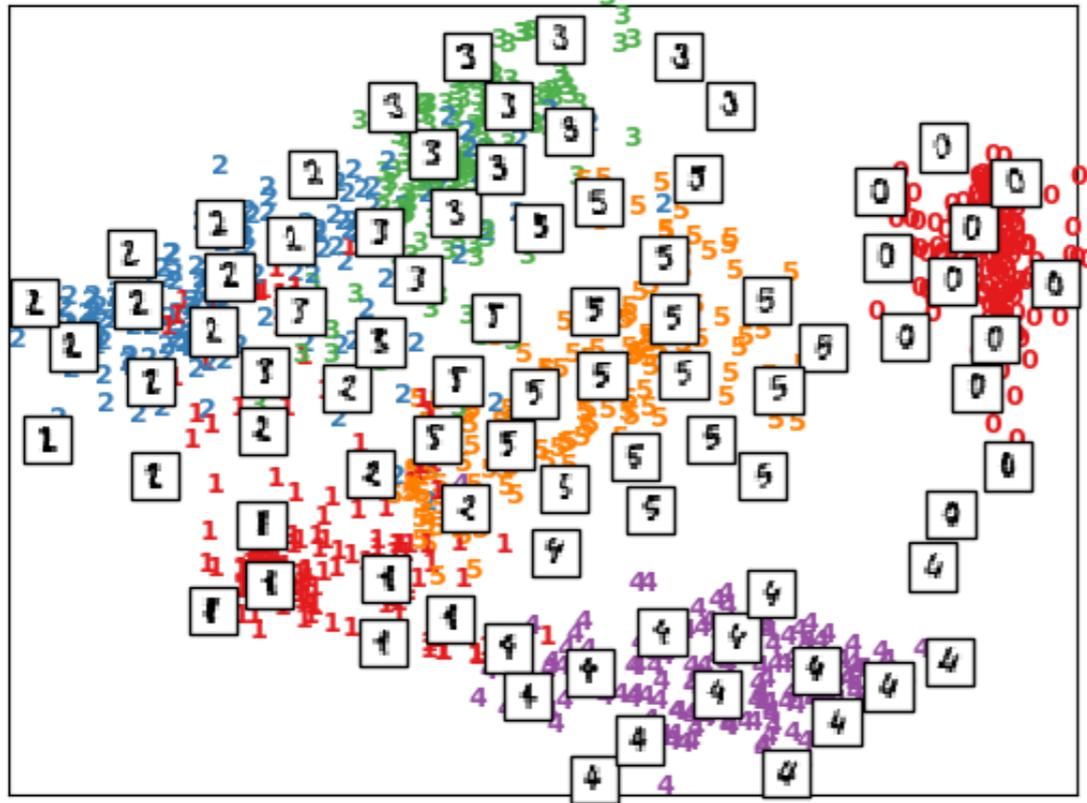
KL-divergence

Manifold Learning

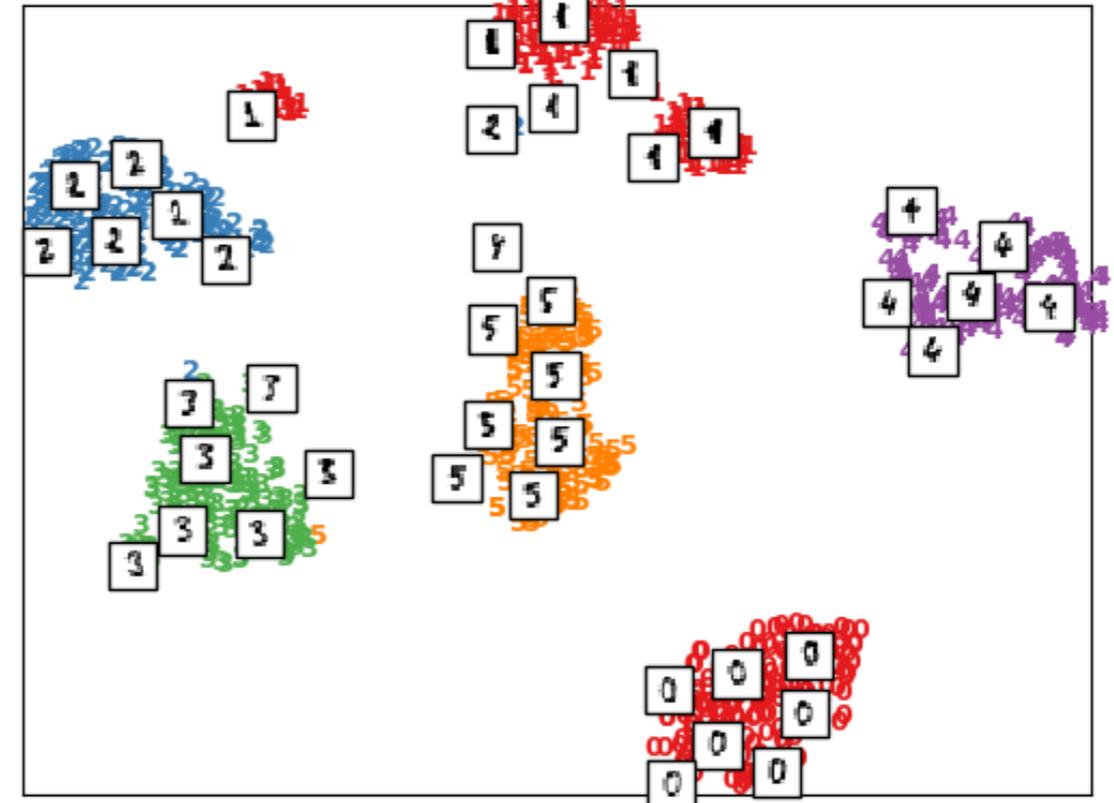
t-SNE chooses a particular generative model for data and for the compression

$$x \in \mathbb{R}^k; y \in \mathbb{R}^l; l \ll k; F(x) \mapsto y$$

Isomap projection of the digits (time 1.12s)



t-SNE embedding of the digits (time 17.58s)



2 parameters: l , perplexity (related to sigmas).