

Monday Lecture

## **Dimensionality Reduction:**

Principal component analysis, Nonnegative matrix factorization, and Independent component analysis

Roman Akhmetshyn, Jennifer Glover, and Teo Lemay

# Dimensionality

# What is a “dimension”

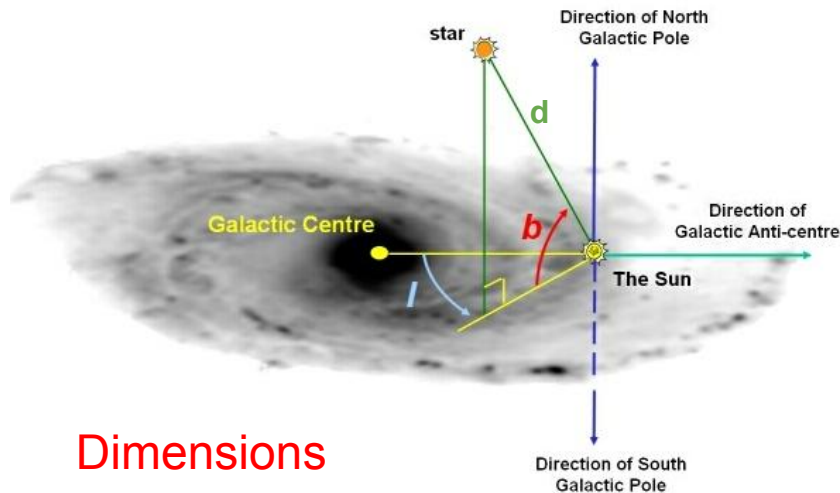
Each of the “things” we can measure about an object is a “**dimension**” or “**feature**” of that object

Let’s imagine we want to know the locations of stars in the galaxy.

Intuitively, we would measure 3 coordinates.

Dimensions

Star	Latitude	Longitude	Distance
1			
2			



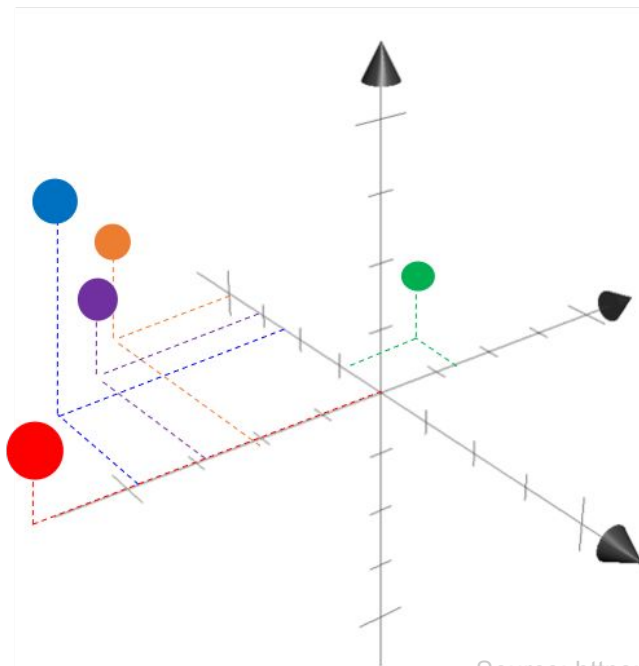
Source: COSMOS

Dimensions

$$\vec{x} = (\ell, b, d)$$

# Dimension of a Space

We need  $D$  pieces of information to describe the location of each object in a  $D$ -dimensional space




$$\vec{x} = (x, y, z)$$

# Intrinsic Dimensionality

- The **intrinsic** dimensionality of a space is the minimum number of dimensions you need to describe a location in the space
- The dimensionality of your data is the number of dimensions you actually measure

Dimensionality of Data = 6

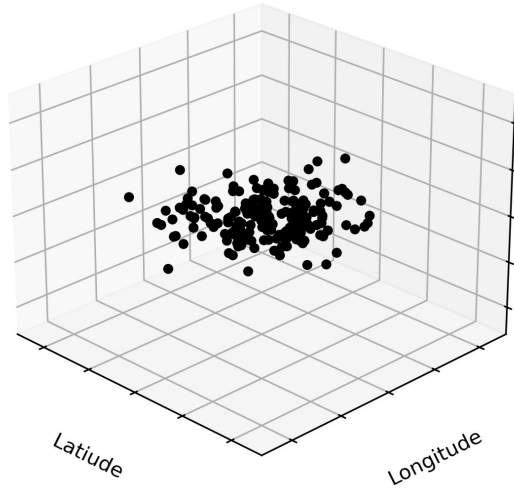


Star	Latitude	Longitude	Distance	RA	DEC	Parallax
1						
2						

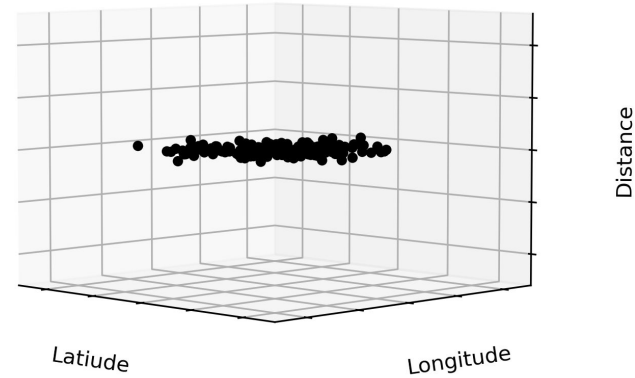
Intrinsic Dimensionality = 3 ?

# What if your data looked like this?

View 1



View 2

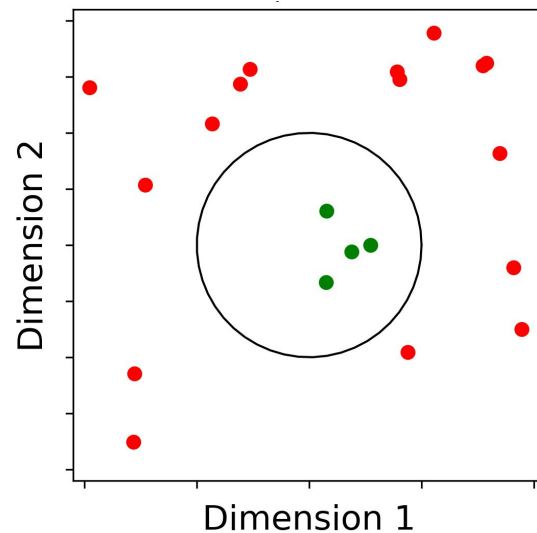
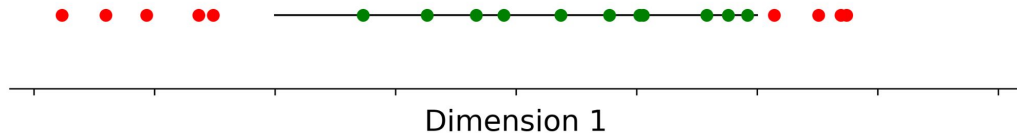


# The “Curse of Dimensionality”

The more dimensions a dataset has, the more data is required to constrain a model

The more dimensions you have:

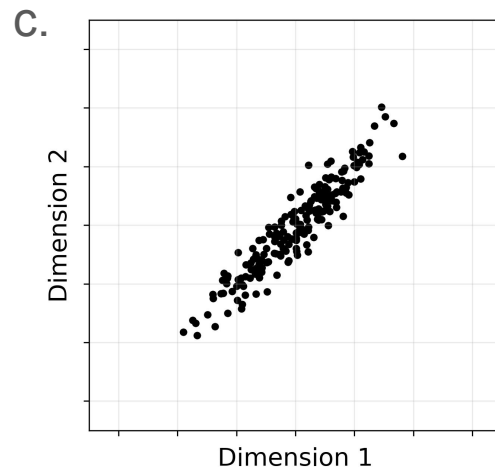
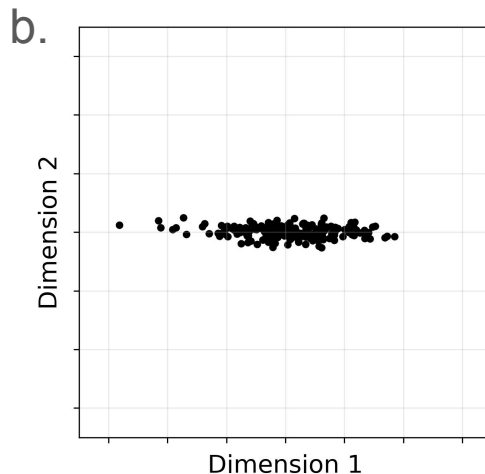
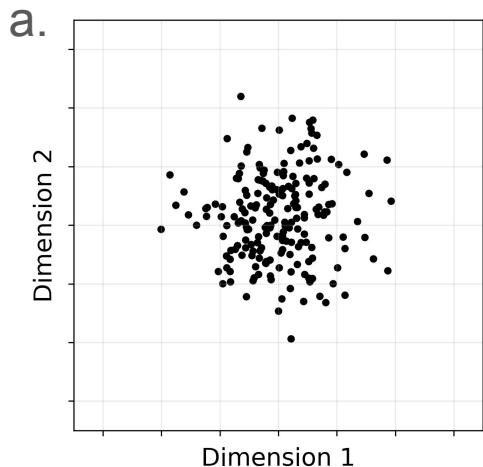
- the more data you need to take
- the more computationally intensive things become
- the more difficult it becomes to visualize the data



# Breaking the “Curse of Dimensionality”

→ Dimension Reduction!

- Ideally, we want to use the minimum required number of dimensions
- There will be some dimensions of the data which capture the most information
  - In real data, many dimensions will be correlated





# Dimension Reduction

# Dimension Reduction

The goal of “**Dimension Reduction**” is to reduce the complexity of the data down to only the most important features

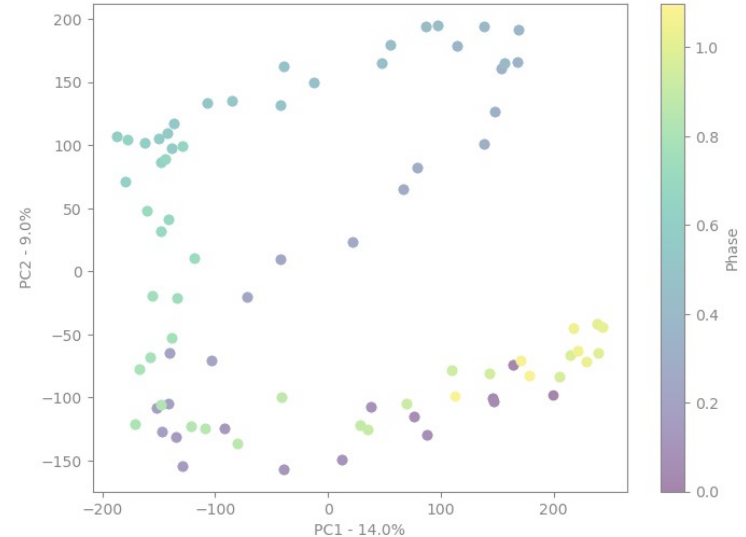
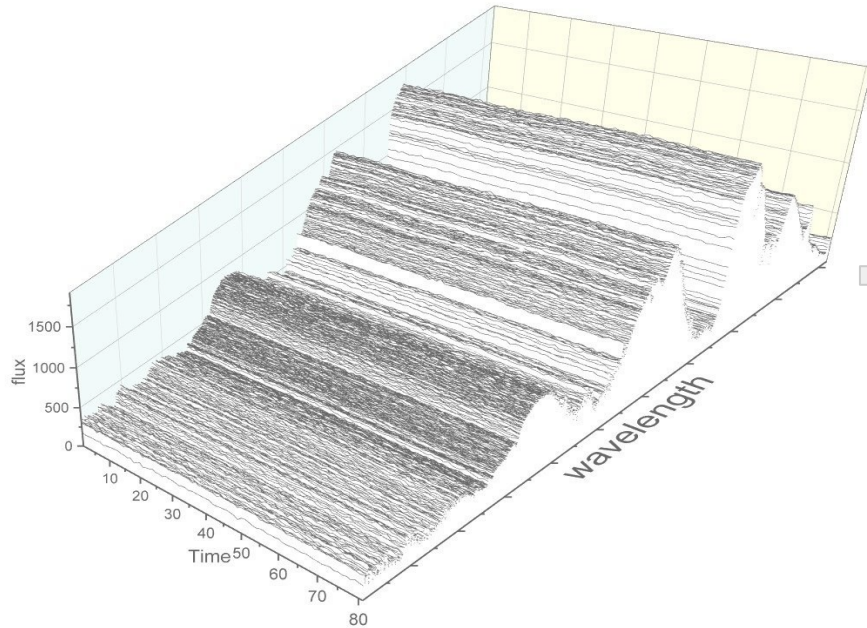
There are several ways to perform dimension reduction

We will go over these three:

- PCA: Principal Component Analysis
- NMF: Nonnegative Component Analysis
- ICA: Independent Component Analysis

Dimension Reduction Method 1:

# Principal Component Analysis (PCA)



# PCA introduction

PCA – linear transformation of multidimensional data that defines orthogonal set of axes which capture most variance.

Generally used to: decorrelate variables, compress data, extract features, and reduce noise.

Important:

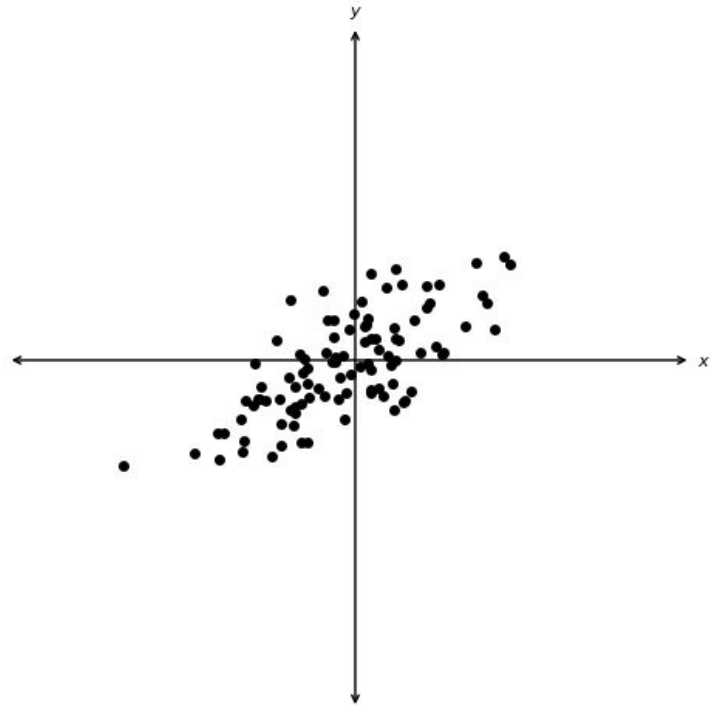
- pre-PCA data must be centered,
- number of PCs =  $\min(\text{measurements}, \text{features})$
- PCs are unphysical
- Data is usually projected on PCs

# How does PCA work?

Before we start, let's center our data by removing the mean.

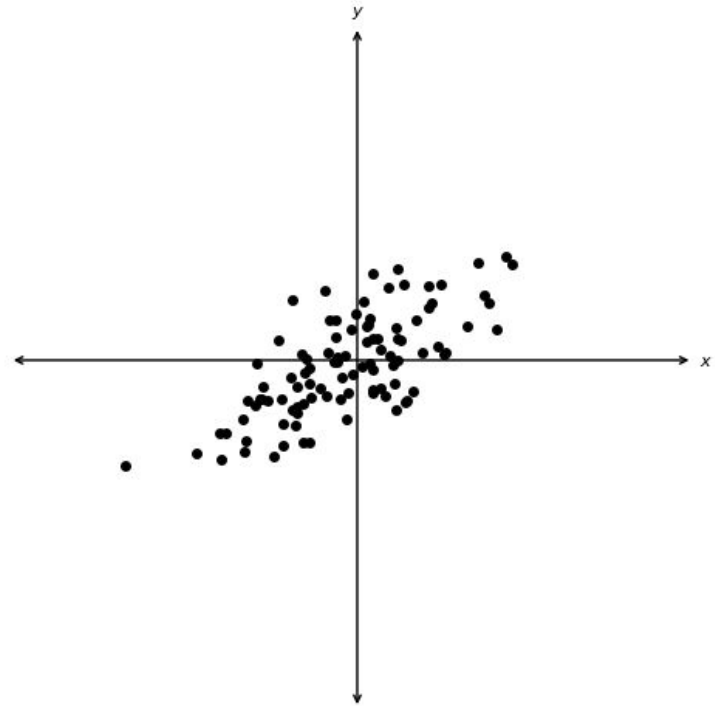
If your dimensions have different units, we should also scale them:

$$X = (\text{data} - \text{mean}) / \text{std}$$



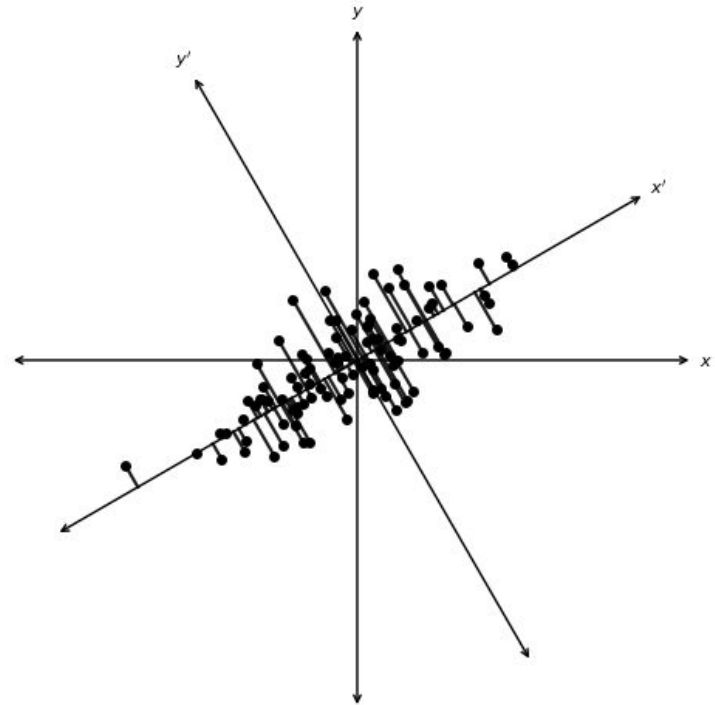
# Information we can infer:

- Correlation of X and Y
- Variance of data
- Clustering
- Patterns



We rotate axes to align with data's correlation, i.e. minimize orthogonal distances of data points to the axes, i.e. maximize variance along the axes.

- Eigenvectors: unit-vectors, our “new” axes, principal components
- Eigenvalues: variance along principal components



# What's next?

PCA is a linear transformation and is based on data projection (orthogonal)

We can:

- Project data onto principal components
- Visualize and explore data with principal component planes
- Reconstruct whole data using few top PCs
- Reconstruct single data using PCP coordinates
- Reconstruct data into less dimensions



# How can we use that in astronomy?

We can:

- Remove low amplitude noise from data (reconstruct by top PCs)
- Remove high amplitude noise (reconstruct by bottom PCs)
- Get a basis to describe objects, standardize parameters
- Optimize models and simulations (reconstruct to lower dimension)
- Prepare data for classification, clustering and ML algorithms
- Interpolate across missing data
- And probably much more!

# PCA pros/cons

## Pros

- Simple
- Computationally efficient
- PCs can be informative
- Handles multicollinearity between variables
- Minimizes reconstruction error

## Cons

- Not scaling or unit invariant
- Need to decide about best number of PCs
- PCs are not always informative
- Fails to capture correlated and non-linear patterns
- Sensitive to outliers
- Not optimal for categorical data

# Monday Exercises:

1. Investigate the “Curse of Dimensionality”
2. Write your own PCA

# PCA Exercises Key steps:

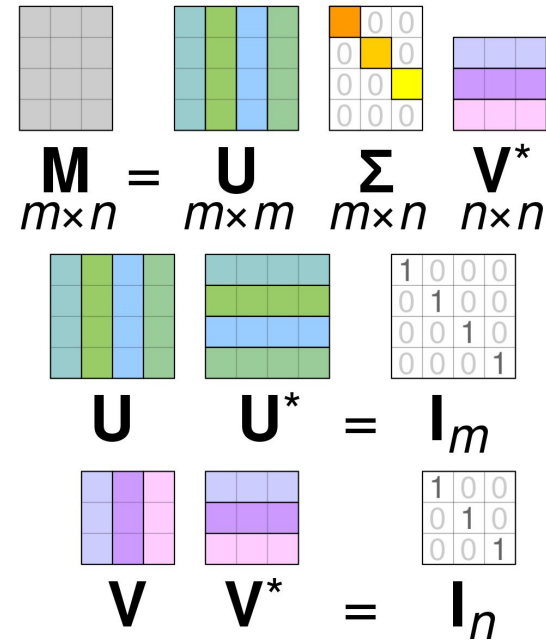
We will use SDSS rest frame spectra of 4000 galaxies that span across 1000 wavelength. Each wavelength is a dimension.

1. We will center data by subtracting mean spectra
2. Use Singular Value Decomposition to retrieve principal components.

# SVD

A method that decomposes a matrix into a product of matrices. A Basis of PSA.

$$U\Sigma V^T = \frac{1}{\sqrt{N-1}}X$$



Source: Wikipedia

**SVD**      $U\Sigma V^T = \frac{1}{\sqrt{N-1}} X$

U - hierarchically arranged basis-columns that describe variance in X-columns

$\Sigma$  - hierarchically arranged non-negative diagonal singular values

V - hierarchically arranged variance-columns of each U-column across data / mixture-columns of U to create X-columns

U and V are unitary ( $UU^T=U^TU=I_{n \times n}$ ,  $VV^T=V^TV=I_{m \times m}$ ), columns are orthonormal.

$$\text{SVD} \quad U \Sigma V^T = \frac{1}{\sqrt{N-1}} X$$

$$X = \underbrace{\begin{bmatrix} | & | & & | \\ x_1 & x_2 & \dots & x_m \\ | & | & & | \end{bmatrix}}_m \Bigg\}^n = U \Sigma V^T = \begin{bmatrix} | & | & & | \\ u_1 & u_2 & \dots & u_n \\ | & | & & | \end{bmatrix}^{n \times n} \begin{bmatrix} \sigma_1 & & & \\ & \sigma_2 & & \\ & & \ddots & \\ & & & \sigma_m \\ \hline & & & & 0 \end{bmatrix}^{n \times m} \begin{bmatrix} - & v_1^T & - \\ - & v_2^T & - \\ & \vdots & \\ - & v_m^T & - \end{bmatrix}^{m \times m}$$

$$= \sigma_1 u_1 v_1^T + \sigma_2 u_2 v_2^T \dots \sigma_m u_m v_m^T$$





SVD  $U\Sigma V^T = \frac{1}{\sqrt{N-1}} X$

$$X = U\Sigma V^T$$

$$X^T X = V \Sigma^T \underbrace{U^T U}_1 \Sigma V^T = V \Sigma^T \Sigma V^T = V \Sigma^2 V^T \quad (\times V)$$

$$\Rightarrow (X^T X)V = V \Sigma^2$$

$$\boxed{X^T X} \quad \begin{array}{c} \text{Ei-vecs} \\ \downarrow \\ \begin{bmatrix} | & | & | \\ v_1 & v_2 & \dots & v_m \\ | & | & | \end{bmatrix} \end{array} = \begin{array}{c} \begin{bmatrix} | & | & | \\ v_1 & v_2 & \dots & v_m \\ | & | & | \end{bmatrix} \end{array} \quad \begin{array}{c} \text{Ei-vals} \\ \downarrow \\ \begin{bmatrix} \sigma^2 & & \\ & \sigma^2 & \\ & & \ddots \\ & & & \sigma_0^2 \end{bmatrix} \end{array}$$

SVD  $U\Sigma V^T = \frac{1}{\sqrt{N-1}} X$

$$X = U\Sigma V^T$$

$$X^T X = V \Sigma^T \underbrace{U^T U}_1 \Sigma V^T = V \Sigma^T \Sigma V^T = V \Sigma^2 V^T \quad (\times V)$$

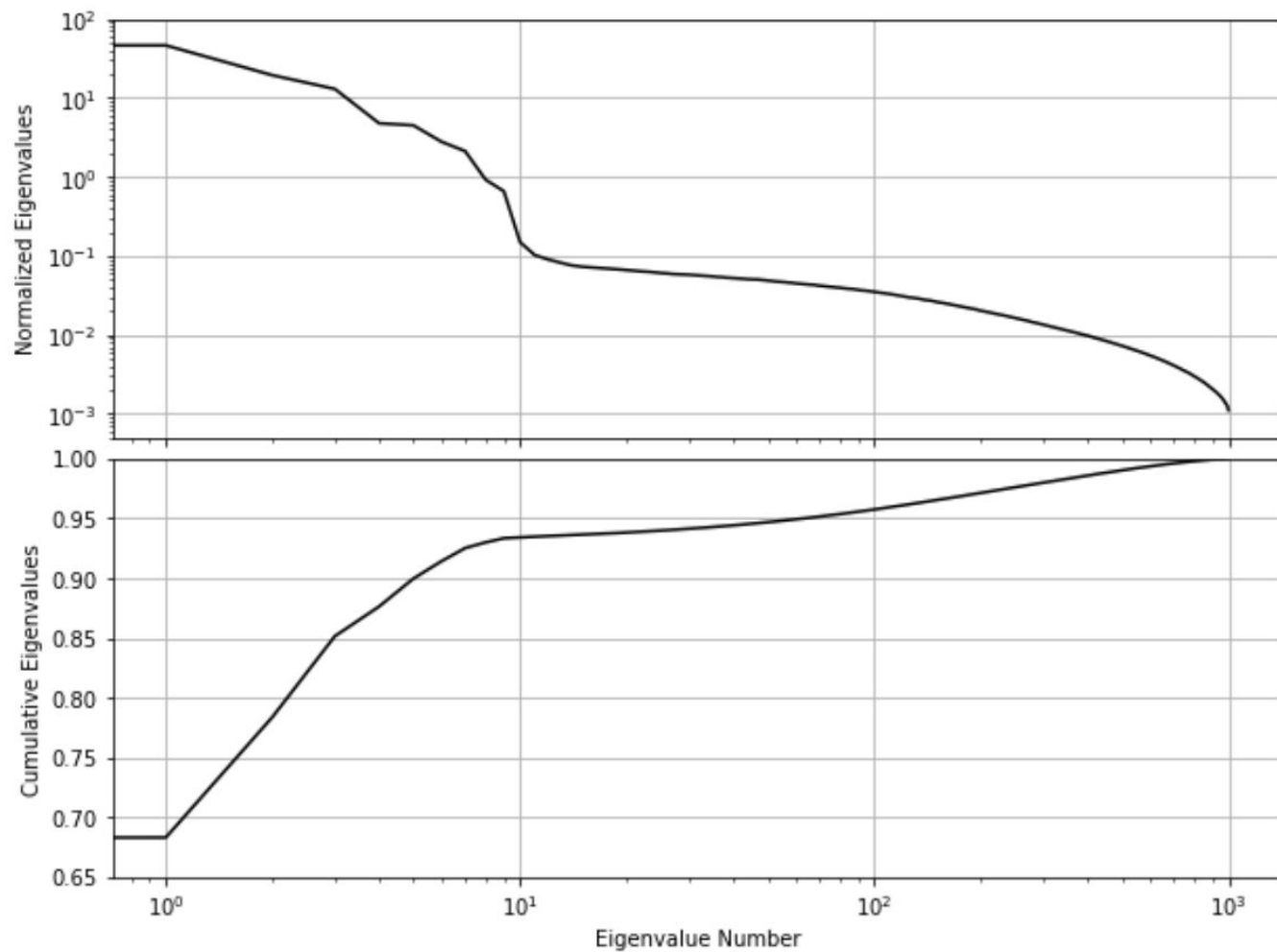
$$\Rightarrow (X^T X)V = V \Sigma^2$$

$$\begin{array}{c} \text{C} \downarrow \\ \boxed{X^T X} \end{array} \quad \begin{array}{c} \text{Ei-vecs} \downarrow \\ \boxed{\begin{array}{c} | \\ v_1 \\ | \\ v_2 \\ | \\ \vdots \\ | \\ v_m \\ | \end{array}} \end{array} = \begin{array}{c} \boxed{\begin{array}{c} | \\ v_1 \\ | \\ v_2 \\ | \\ \vdots \\ | \\ v_m \\ | \end{array}} \end{array} \quad \begin{array}{c} \text{Ei-vals} \downarrow \\ \boxed{\begin{array}{c} \sigma^2 \\ \sigma^2 \\ \vdots \\ \sigma_0^2 \end{array}} \end{array}$$

# PCA Exercises Key steps:

We will use SDSS rest frame spectra of 4000 galaxies that span across 1000 wavelength. Each wavelength is a dimension.

1. We will center data by subtracting mean spectra
2. Use Single Value Decomposition to retrieve principal components.
3. Scree plot



# PCA Exercises Key steps:

We will use SDSS rest frame spectra of 4000 galaxies that span across 1000 wavelength. Each wavelength is a dimension.

1. We will center data by subtracting mean spectra
2. Use Single Value Decomposition to retrieve principal components.
3. Scree plot
4. Project data onto PCP
5. Reconstruct data

**Dimensionality reduction:**

Principal component analysis, Nonnegative matrix factorization, and Independent component analysis

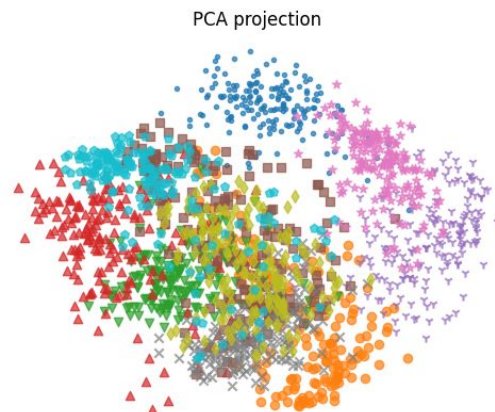
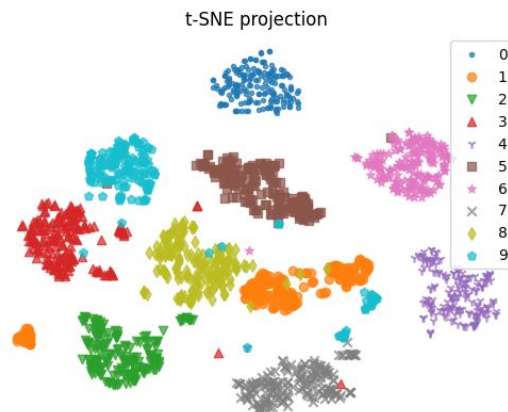
Roman Akhmetshyn, Jennifer Glover, and Teo Lemay

# PCA is not the only option

- Nonnegative Matrix Factorization (NMF)
- Independent component analysis (ICA)
- Manifold learning
- And more! (Kernel methods, Autoencoder, t-SNE, etc.)

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



# “Important Component” methods

- **Nonnegative Matrix Factorization (NMF)**
- **Independent component analysis (ICA)**
- Manifold learning
- And more! (Kernel methods, Autoencoder, t-SNE, etc.)

$$\text{Data} = \mathbf{X} \approx \mathbf{WH}$$

$$\mathbf{W} \rightarrow \text{Components}$$

$$\mathbf{H} \rightsquigarrow \text{Weights}$$



Dimension Reduction Method 2:

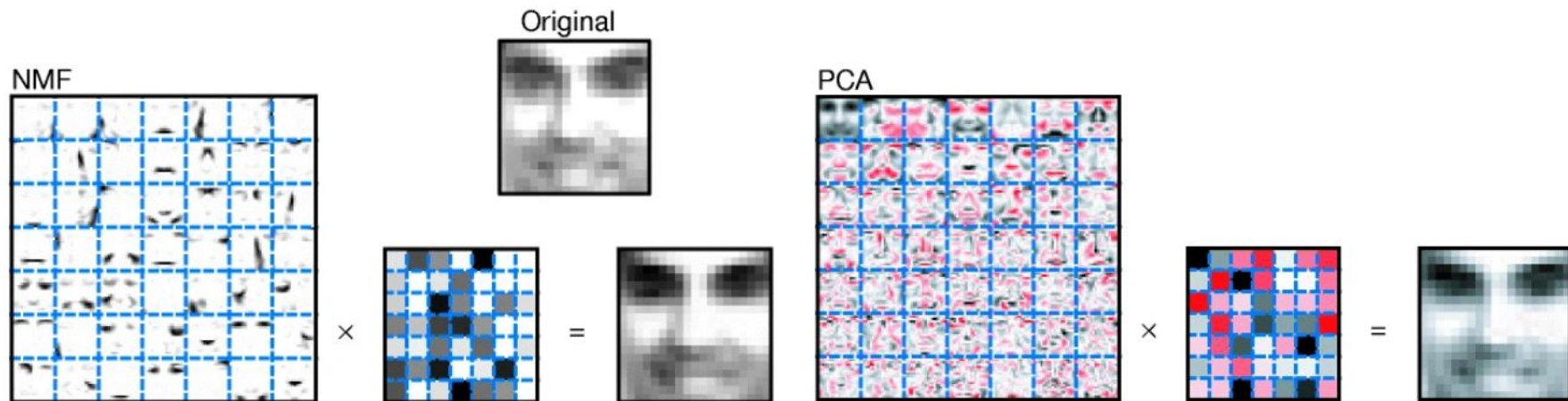
# Nonnegative Matrix Factorization (NMF)

# Nonnegative Matrix Factorization (NMF)

- Very similar to PCA
  - Dimensionality reduction
  - Sometimes informative components
- Constrain solution with known characteristics of the data
- A lot of astronomical data is non-negative ...

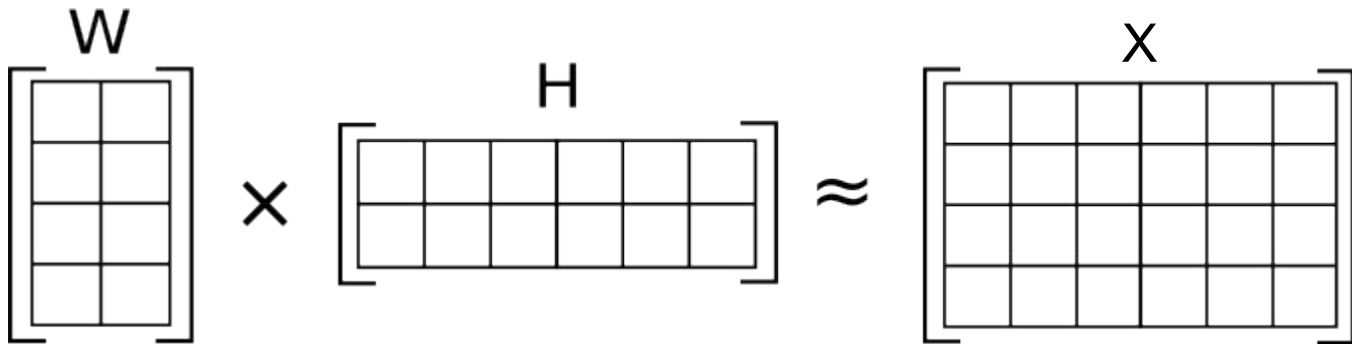
# Nonnegative Matrix Factorization (NMF)

- NMF constraints mean that different components get selected
- PCA tends towards holistic features, NMF can pick out important parts



# How does NMF work?

- Choose number of intrinsic dimensions  $<$  actual dimensionality
- Guess initial values, calculate distance metric



$$\begin{aligned}\chi^2 &= \|(\mathbf{X} - \mathbf{W}\mathbf{H})\|^2 \\ &= \sum_{ij} \left( X_{ij} - \sum_k W_{ik} H_{kj} \right)^2\end{aligned}$$

# How does NMF work?

- Choose number of intrinsic dimensions  $<$  actual dimensionality
- Guess initial values, calculate distance metric
- **Iterative update rule**

$$\begin{aligned}\chi^2 &= \|(\mathbf{X} - \mathbf{W}\mathbf{H})\|^2 \\ &= \sum_{ij} \left( X_{ij} - \sum_k W_{ik} H_{kj} \right)^2\end{aligned}$$

$$\begin{aligned}\mathbf{H} &\leftarrow \mathbf{H} \circ \frac{\mathbf{W}^T \mathbf{X}}{\mathbf{W}^T \mathbf{W} \mathbf{H}} \\ \mathbf{W} &\leftarrow \mathbf{W} \circ \frac{\mathbf{X} \mathbf{H}^T}{\mathbf{W} \mathbf{H} \mathbf{H}^T}\end{aligned}$$

# How does NMF work?

- Choose number of intrinsic dimensions < actual dimensionality
- Guess initial values, calculate distance metric
- Iterative update rule
- **NMF is sensitive to local minima**
- **Multiple random initializations**

$$\begin{aligned}\chi^2 &= \|(\mathbf{X} - \mathbf{W}\mathbf{H})\|^2 \\ &= \sum_{ij} \left( X_{ij} - \sum_k W_{ik} H_{kj} \right)^2\end{aligned}$$

$$\begin{aligned}\mathbf{H} &\leftarrow \mathbf{H} \circ \frac{\mathbf{W}^T \mathbf{X}}{\mathbf{W}^T \mathbf{W} \mathbf{H}} \\ \mathbf{W} &\leftarrow \mathbf{W} \circ \frac{\mathbf{X} \mathbf{H}^T}{\mathbf{W} \mathbf{H} \mathbf{H}^T}\end{aligned}$$

# How does NMF work?

- Choose number of intrinsic dimensions  $<$  actual dimensionality
- Guess initial values, calculate distance metric
- Iterative update rule
- **NMF is sensitive to local minima**
- **Multiple random initializations**

$$\chi^2 = \sum_{i,j} \| \mathbf{x}_{ij} - \sum_k \mathbf{w}_k \mathbf{h}_k \|^2$$

`sklearn.decomposition.NMF`

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

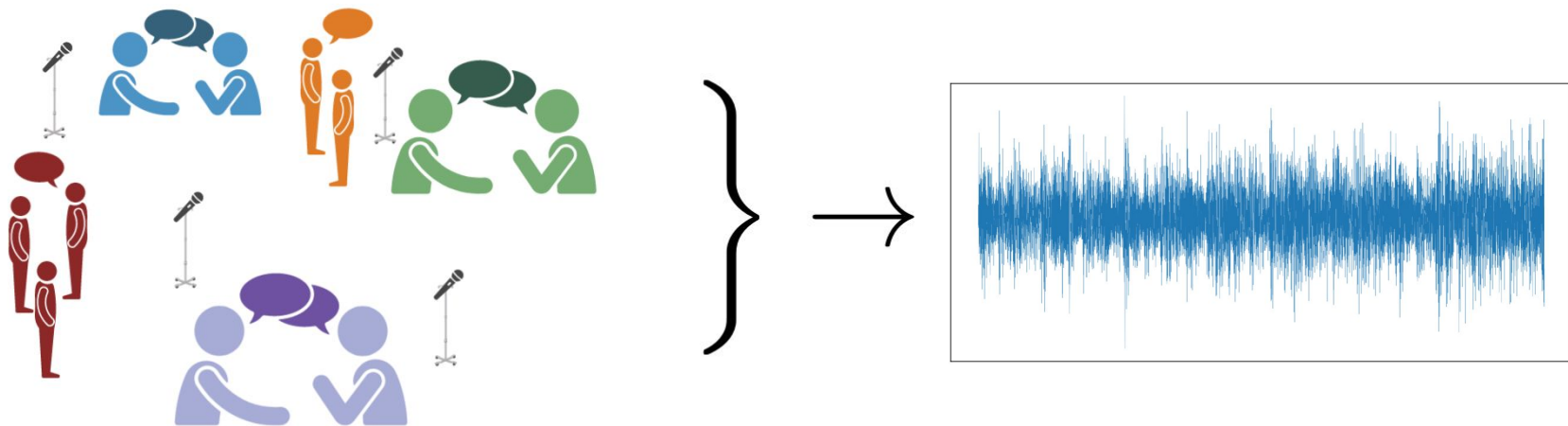
Dimension Reduction Method 3:

# Independent Component Analysis (ICA)



# Independent Component Analysis (ICA)

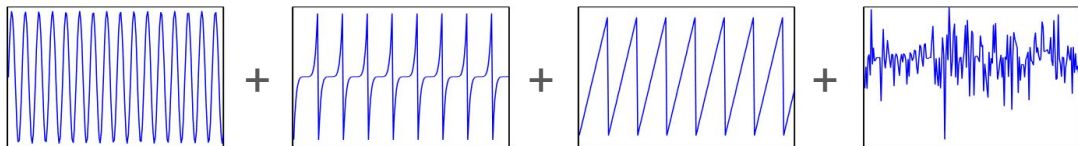
- Another matrix decomposition approach
  - Emphasis on important components
  - Exploration of patterns in data
- Assume data is a combination of signals from independent sources



# Independent Component Analysis (ICA)

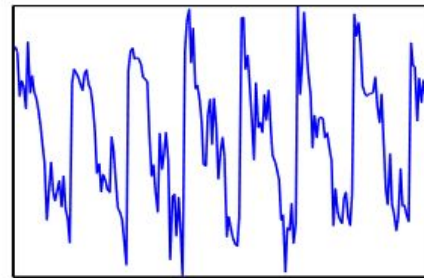
## ICA

- Find statistically independent “source signals”
- Physically interpretable Independent Components?

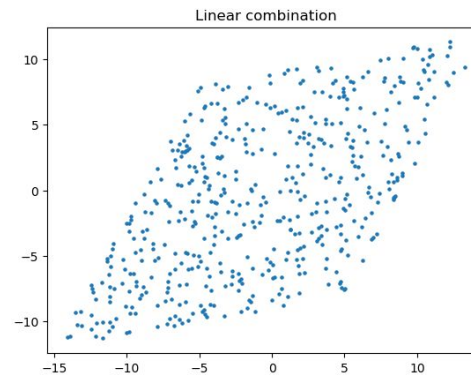
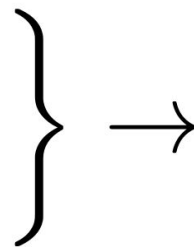
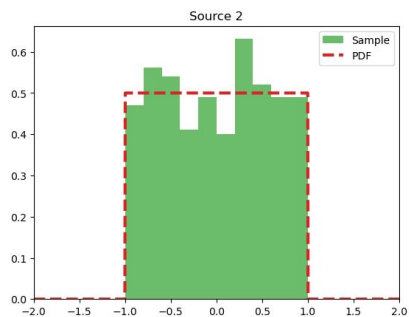
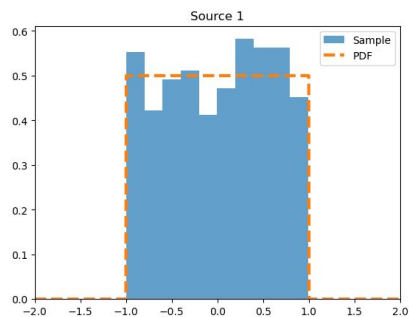


## PCA

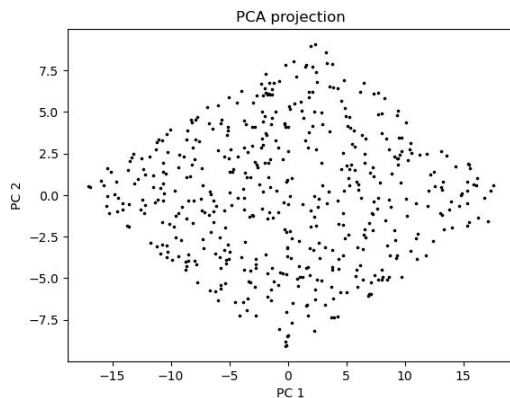
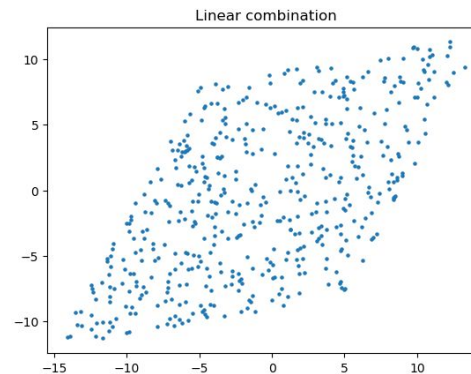
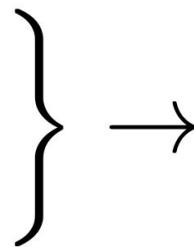
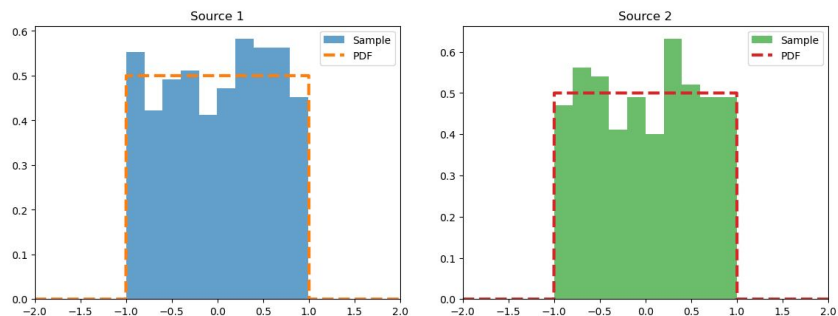
- Maximize variance along principle components
- Eigenvectors of covariance matrix for the data



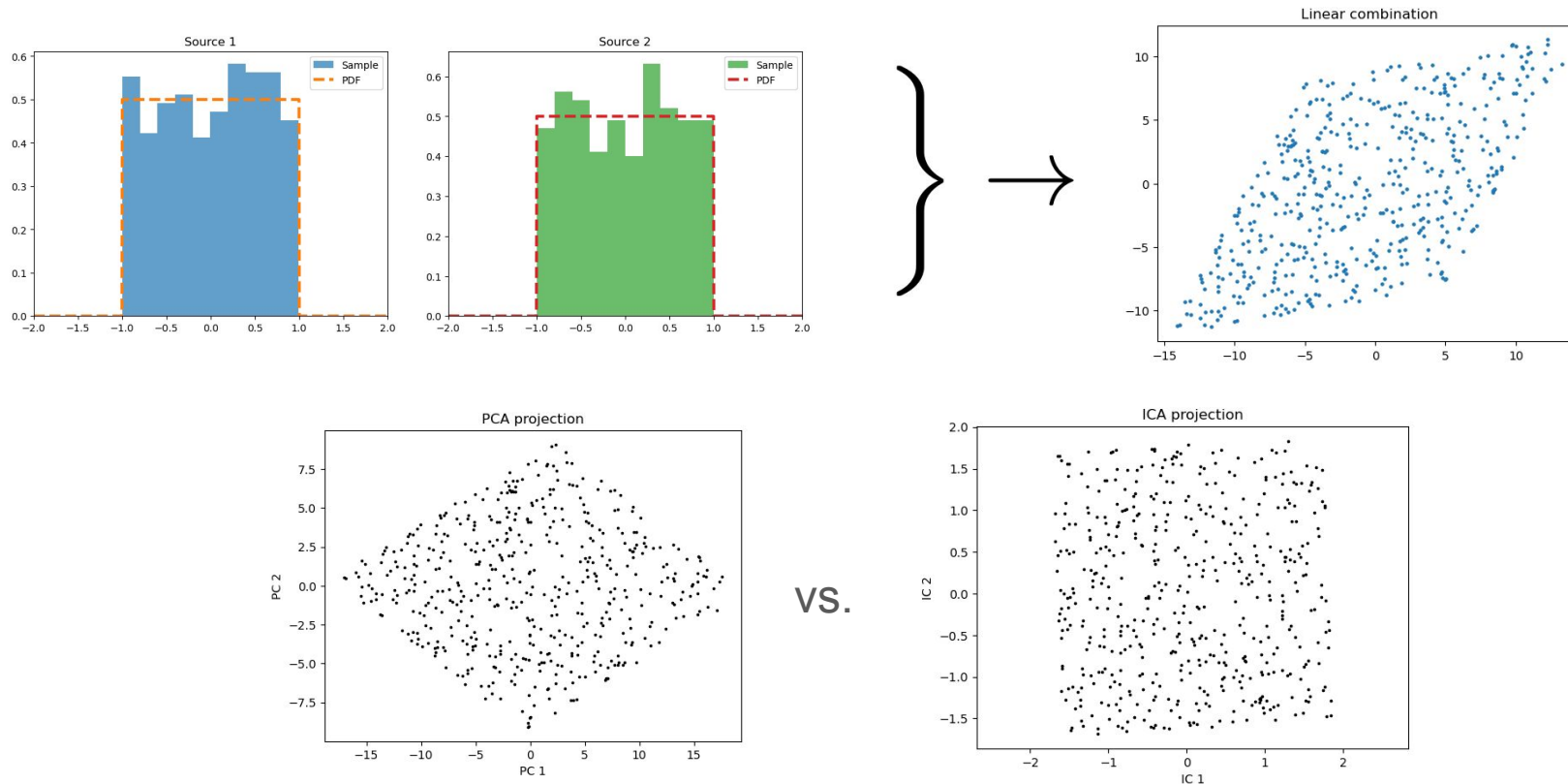
# Independent Component Analysis (ICA)



# Independent Component Analysis (ICA)



# Independent Component Analysis (ICA)



# Independent Component Analysis (ICA)

## Assumptions

- Sources in the data are independent
- Sources are non-gaussian

$$p(y_a, y_b, y_c, \dots) = p(y_a)p(y_b)p(y_c) \dots$$

## Notes

- No measure of variance explained by each component
- No hierarchy of independent components
- Should center data (mean subtraction)

# How does ICA work?

- Central limit theorem!
- Combinations of random variables tends towards a gaussian
- Assume data should be more gaussian than the sources
- Choose the number Independent components to look for
- Try to find the maximally **non**-gaussian matrix decomposition

$$\mathbf{X} = \mathbf{WH}$$

*PCA*    $\mathbf{W}^{-1}\mathbf{X}$    maximizes variance

*ICA*    $\mathbf{W}^{-1}\mathbf{X}$    maximizes non-gaussianness

# Independent Component Analysis (ICA)

- A few different measures of non-gaussian-ness
  - Kurtosis, higher power cumulants
  - Neg-entropy
  - Mutual information minimization
- Practical implementations (fastICA) combine parts of different measures and make useful approximations
- Requires setting the number of ICs to look for



# Independent Component Analysis (ICA)

- A few different measures of non-gaussian-ness
  - Kurtosis, higher power cumulants
  - Neg-entropy
  - Mutual information minimization
- Practical implementations (fastICA) combine parts of different measures and make useful approximations
- Requires setting the number of ICs to look for

`sklearn.decomposition.FastICA`

# PCA is not the only option!

- PCA is the simplest, and most widely applicable of these methods
- What do you know about your data a priori?
  - Noisy combination of signals?
  - Nonlinear high dimensional shape?
- What are you trying to achieve with dimensionality reduction?
  - Do the components matter?
  - Should you believe the components?

# Wednesday Exercise: Play around with parameters

- NMF example with Olivetti faces dataset
- ICA mixtures of waveforms
- Find sources in a mystery dataset