# Dimensionality Reduction

By Alex Cayco Gajic

neuromatch academy

# Who is Alex Cayco Gajic?

- Junior Professor @ ENS (Paris)

- Motor control, cerebellar theory
- Population coding, statistical learning

# Overview of tutorials

1. Geometric view of data
2. Principal component analysis

**Toy data (2D)**

1. Dimensionality reduction and reconstruction
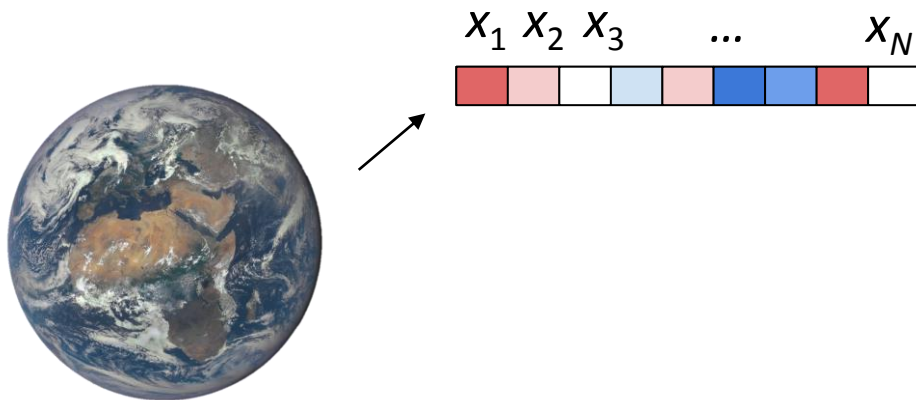2. Nonlinear dimensionality reduction

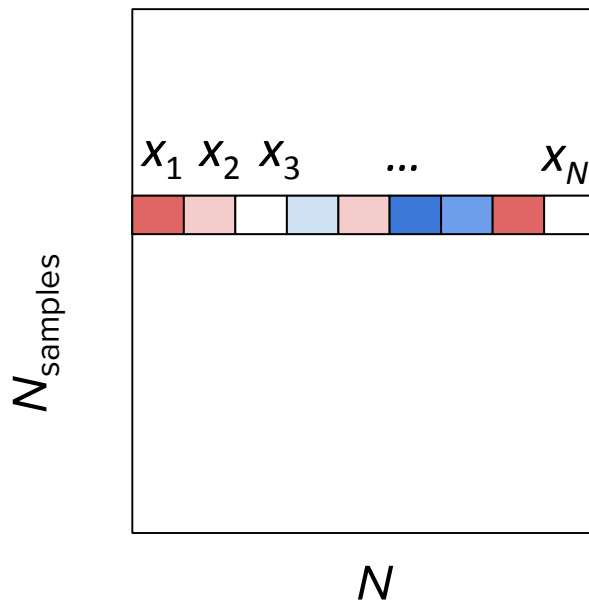**Real data (high-D)**

# Geometric view of data

Tutorial 1

# Multivariate data

$x_1$ $x_2$ $x_3$ ... $x_N$

Sample of $N$ variables during a single observation.

# How to represent multivariate data?

$$\mathbf{X} =$$



$x_1$ $x_2$ $x_3$ ... $x_N$

$N_\text{samples}$

$N$

Single sample of all variables

# How to represent multivariate data?
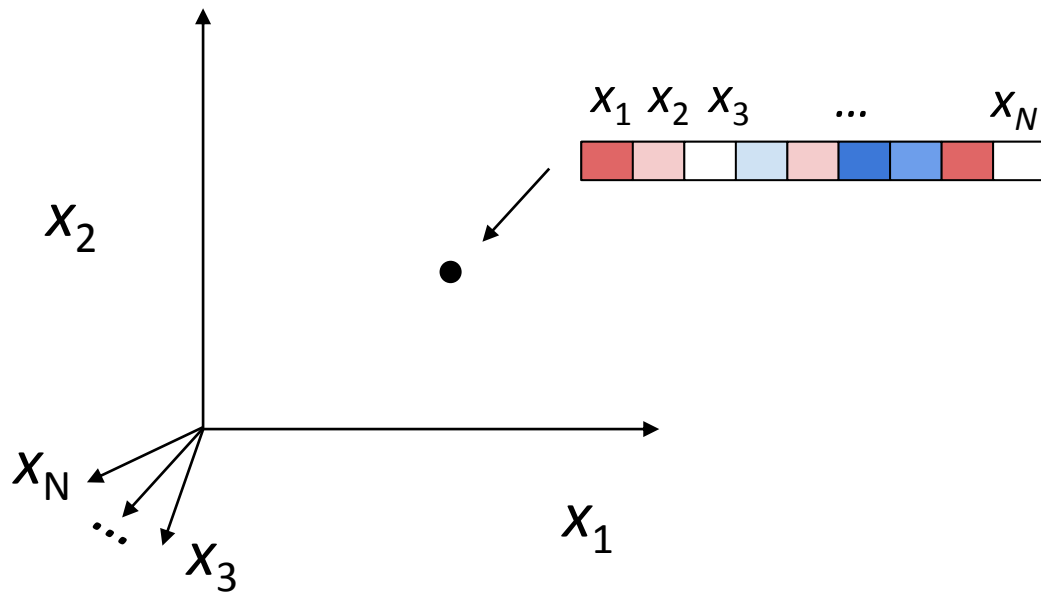
$$\mathbf{X} =$$



$N_{\text{samples}}$

$\boldsymbol{x}_3$

$N$
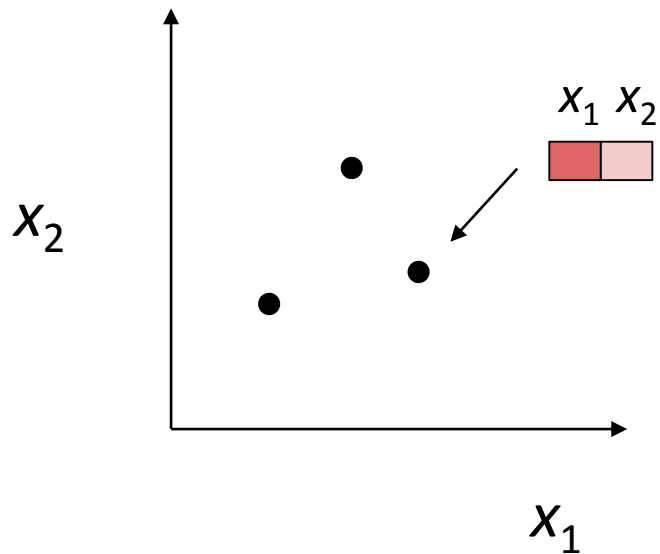
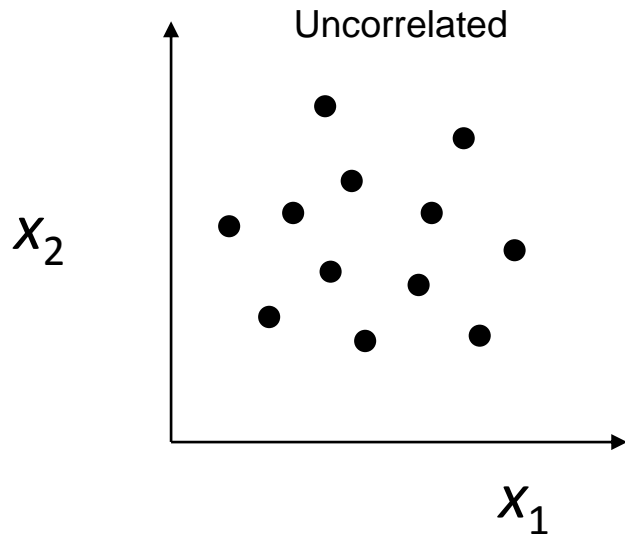**All** samples of a **single** variable

# How to represent multivariate data?

$$x_1 \; x_2 \; x_3 \qquad \ldots \qquad x_N$$

$x_2$

$x_N$

$\ldots$

$x_3$

$x_1$
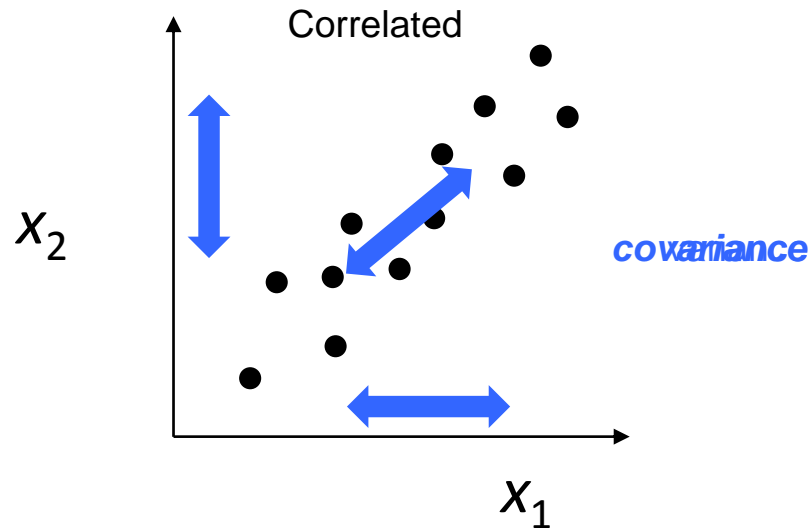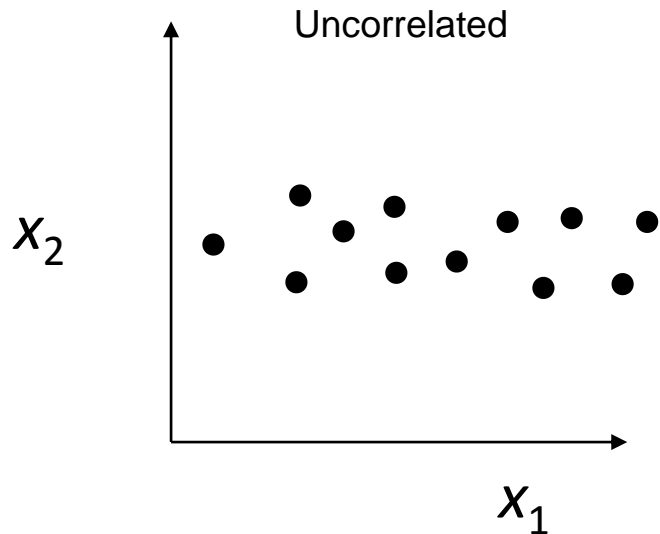
# How to represent multivariate data?

# Multivariate data has variability

Uncorrelated

$x_2$

$x_1$

# Multivariate data has variability



Uncorrelated

$x_2$

$x_1$

Correlated

$x_2$

$x_1$

covariance

# PCA: the big picture

Look for directions of maximum variance.



$x_2$

$x_3$

$x_1$

Component 2

Component 1

# Multivariate data has variability



Uncorrelated

$x_2$

$x_1$

Correlated

$x_2$

$x_1$

# How to quantify variability

**Variance**    $$\mathrm{var}(x_1) = E[x_1^2] - E[x_1]^2$$

**Covariance**    $$\mathrm{cov}(x_1, x_2) = E[x_1 x_2] - E[x_1]E[x_2]$$

**Correlation**    $$\rho = \frac{\mathrm{cov}(x_1, x_2)}{\sqrt{\mathrm{var}(x_1)\mathrm{var}(x_2)}}$$

**Normalized to be within -1 to +1**

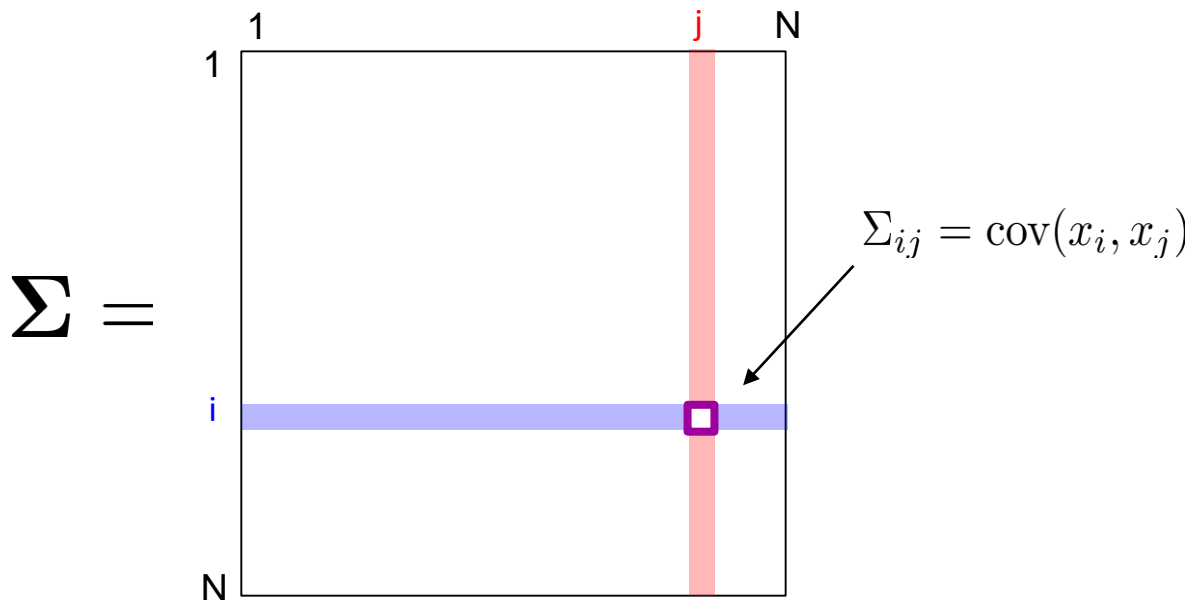# The covariance matrix $\mathbf{\Sigma}$

$$\mathbf{\Sigma} =$$

$$\Sigma_{ij} = \text{cov}(x_i, x_j)$$

# The covariance matrix $\boldsymbol{\Sigma}$



$$\boldsymbol{\Sigma} =$$

$$\Sigma_{ij} = \mathrm{cov}(x_i, x_j)$$

$$\Sigma_{ji} = \mathrm{cov}(x_j, x_i)$$

# The covariance matrix $\Sigma$

$$\Sigma =$$



$$\Sigma_{ii} = \text{var}(x_i)$$

$$\Sigma_{ij} = \text{cov}(x_i, x_j)$$

$$\Sigma_{ji} = \text{cov}(x_j, x_i)$$

# The covariance matrix $\mathbf{\Sigma}$



$$\mathbf{\Sigma} =$$

- Variances on the diagonal
- Covariances on the off-diagonal
- Symmetric matrix

$$\Sigma_{ij} = \Sigma_{ji}$$

$$\mathbf{\Sigma}^T = \mathbf{\Sigma}$$
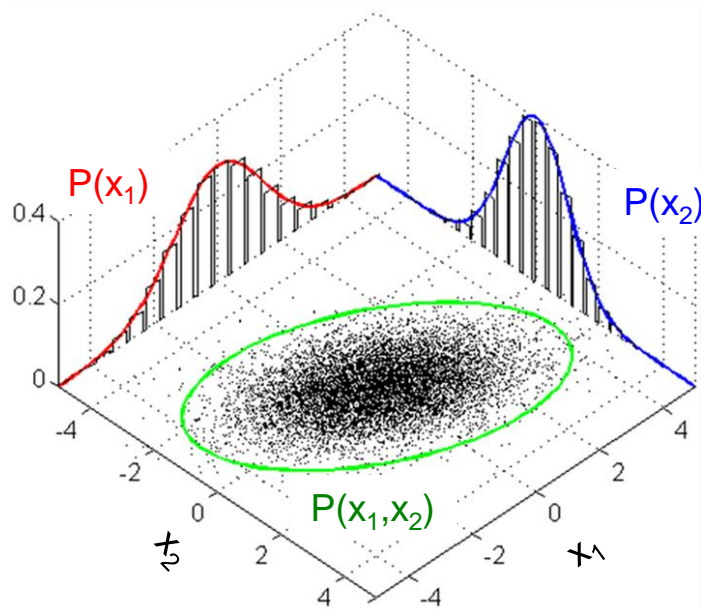
# Multivariate normal distribution

Generalization of normal distribution to N dimensions

$$P(\mathbf{x}) \sim e^{-\frac{1}{2}(\mathbf{x}-\boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x}-\boldsymbol{\mu})}$$

- Parameters:
  $\boldsymbol{\mu}$ - mean of each variable
  $\boldsymbol{\Sigma}$ - covariance matrix
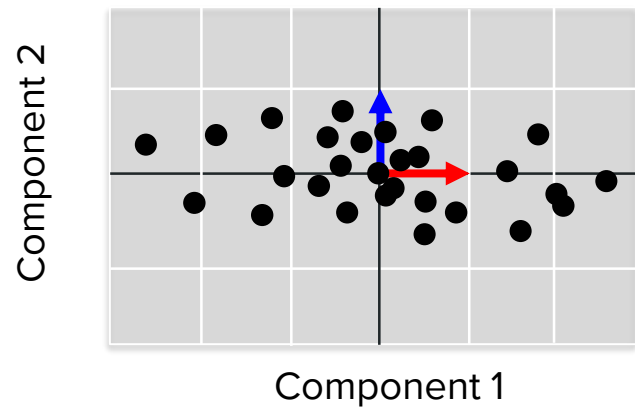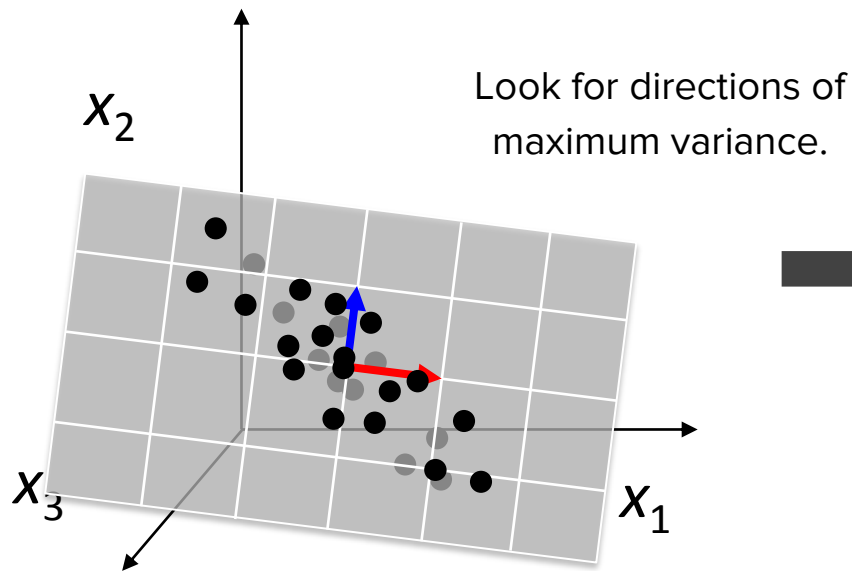
- Marginal distribution $P(x_i)$ is 1D Gaussian

# (break for tutorial exercise)

neuromatch
academy

# PCA: the big picture

Look for directions of maximum variance.

$x_2$

$x_3$

$x_1$

Component 2

Component 1

# Many ways to represent multivariate data

$$w = [0, 2]$$
$$w = 3u + 2w$$

**Basis:**
- Set of N vectors with which you can construct any point in the N-dimensional vector space.

$x_2$

$w$

$u$

$x_1$

# Many ways to represent multivariate data

**v** = [*3.5*,*-1.2*]

x₂ → $x_2$

x₁ → $x_1$

**w**   **u**

**Basis:**
- Set of N vectors with which you can construct any point in the N-dimensional vector space.

**Orthogonal basis**
- All basis vectors are orthogonal.

# Many ways to represent multivariate data

Not orthogonal

$x_2$

$v = [7,-3.8]$

$x_1$

$w$   $u$

**Basis:**
- Set of N vectors with which you can construct any point in the N-dimensional vector space.

**Orthogonal basis**
- All basis vectors are orthogonal.

# Many ways to represent multivariate data

Orthogonal
not orthonormal

$x_2$

$\mathbf{v} = [\textcolor{red}{1},\textcolor{blue}{1}]$

$x_1$

$\mathbf{w}$

$\mathbf{u}$

**Basis:**
- Set of N vectors with which you can construct any point in the N-dimensional vector space.

**Orthogonal basis**
- All basis vectors are orthogonal.
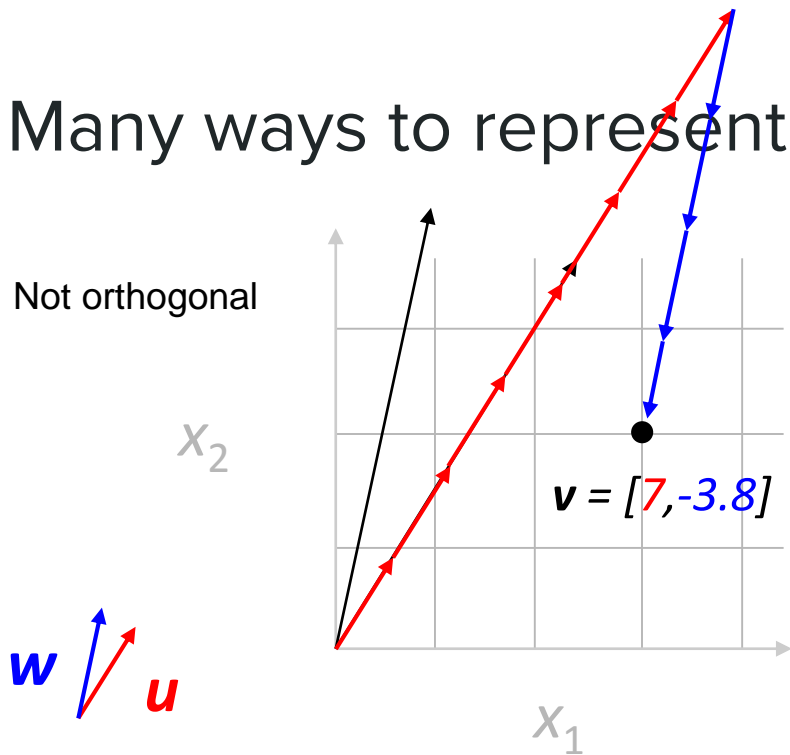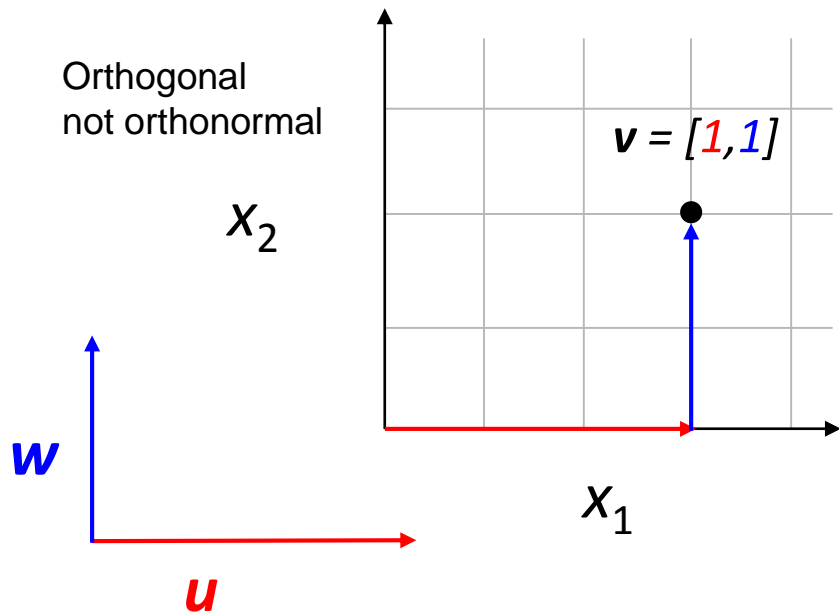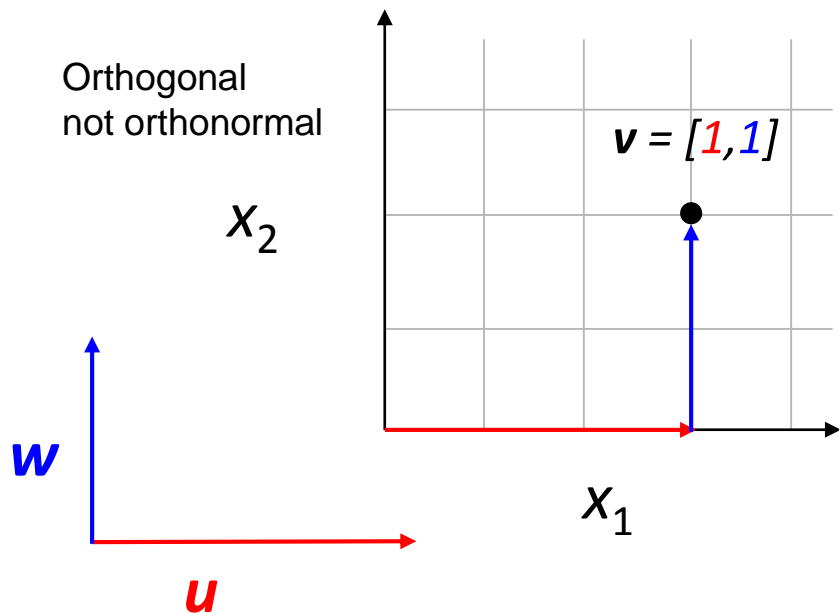
**Orthonormal basis**
- Orthogonal + all basis vectors have a length of 1.

$$\|\mathbf{u}\| = \|\mathbf{w}\| = 1$$

$$\|\mathbf{u}\| = \sqrt{u_1^2 + u_2^2}$$

# Many ways to represent multivariate data

Orthogonal
not orthonormal

$x_2$

$\boldsymbol{v} = [\textcolor{red}{1},\textcolor{blue}{1}]$

$x_1$

**w**

**u**

**Basis:**
- Set of N vectors with which you can construct any point in the N-dimensional vector space.

**Orthogonal basis**
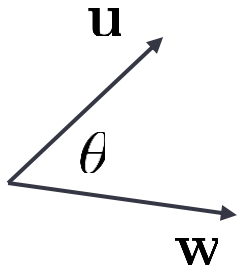- All basis vectors are orthogonal.

**Orthonormal basis**
- Orthogonal + all basis vectors have a length of 1.

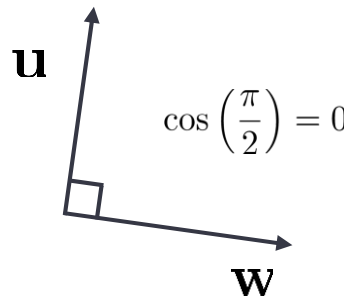An orthogonal basis can easily be **normalized**:

$$\tilde{\mathbf{u}} = \frac{\mathbf{u}}{\|\mathbf{u}\|} \qquad \tilde{\mathbf{w}} = \frac{\mathbf{w}}{\|\mathbf{w}\|}$$

# The dot product

$$\mathbf{u} \cdot \mathbf{w} = \|\mathbf{u}\| \, \|\mathbf{w}\| \cos(\theta)$$



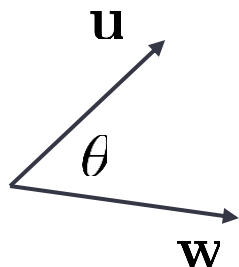$$\cos\left(\frac{\pi}{2}\right) = 0$$

Orthogonal vectors ⬌ Dot product is zero

# The dot product

$$\mathbf{u} \cdot \mathbf{w} = \|\mathbf{u}\| \, \|\mathbf{w}\| \cos(\theta)$$

$$\mathbf{u} \cdot \mathbf{w} = \sum_{i=1}^{N} u_i w_i$$



$$\mathbf{u} \cdot \mathbf{w} \qquad \mathbf{u}^T \qquad \mathbf{w}$$

$$\square \;=\; \boxed{\phantom{xxxxx}}$$

Orthogonal vectors $\longleftrightarrow$ Dot product is zero

$$\mathbf{u}^T \mathbf{w} = 0$$

# (break for tutorial exercise)

# Change of basis



v = ?

[3,2] → [3.5,-1.2]

Standard basis    New basis

How do we transform coordinates to a new orthonormal basis?

# Change of basis via projection

Project **v** onto **u**

$\|\mathbf{v}\| \cos(\theta)$



$\mathbf{u} \cdot \mathbf{v} = \|\mathbf{u}\|\|\mathbf{v}\| \cos(\theta)$

$= \|\mathbf{v}\| \cos(\theta)$

$= \mathbf{u}^T \mathbf{v} = \mathbf{v}^T \mathbf{u}$

New coordinates from dot product!

# Projection to orthonormal basis

For an orthonormal basis,
new coordinates are

$$y_1 \quad = \quad \boxed{x_1 \;\; x_2} \boxed{\textcolor{red}{u}}$$

old
coordinates $\bullet$ new basis
vector

$$y_2 \quad = \quad \boxed{x_1 \;\; x_2} \boxed{\textcolor{blue}{w}}$$

# Projection to orthonormal basis

For an orthonormal basis,
new coordinates are

$$\boxed{y_1 \; y_2} \; = \; \boxed{x_1 \; x_2} \; \boxed{\textcolor{red}{u} \; \textcolor{blue}{w}}$$

$$\underset{\substack{\text{old} \\ \text{coordinates}}}{} \bullet \underset{\substack{\text{new basis} \\ \text{vector}}}{}$$

# Projection to orthonormal basis

$$\mathbf{Y} = \mathbf{X} \quad \mathbf{W}$$

For an orthonormal basis, new coordinates are

$$\text{old coordinates} \cdot \text{new basis vector}$$

$$\begin{array}{|cc|}\hline y_1 & y_2 \\\hline \\\hline \\\hline \end{array} = \begin{array}{|cc|}\hline x_1 & x_2 \\\hline \\\hline \\\hline \end{array}\begin{array}{|c|c|}\hline u & w \\\hline\end{array}$$

adding more samples…

# Projection to orthonormal basis

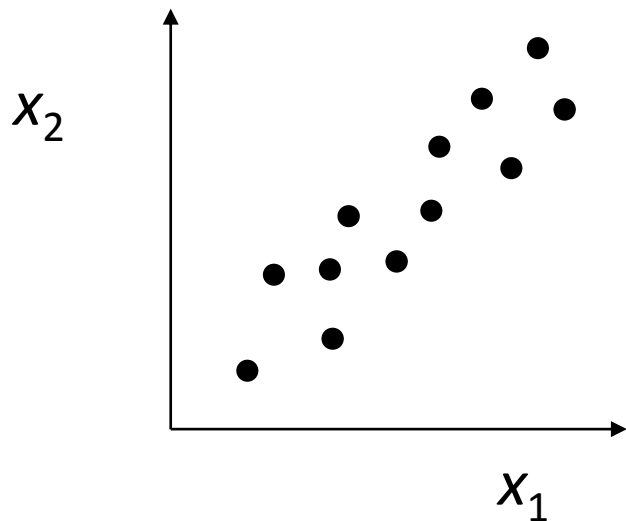$$\mathbf{Y} = \mathbf{X}\,\mathbf{W}$$

# Principal components analysis

Tutorial 2

# Goal of dimensionality reduction

*N variables*                    *K features*

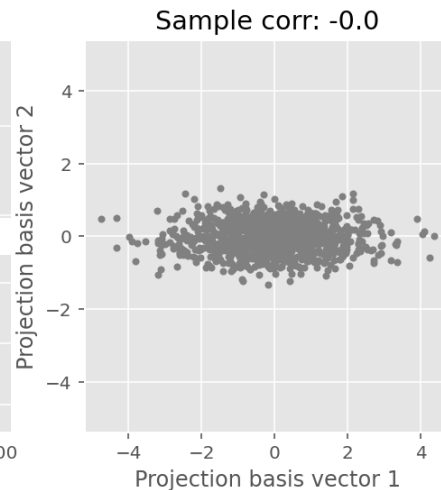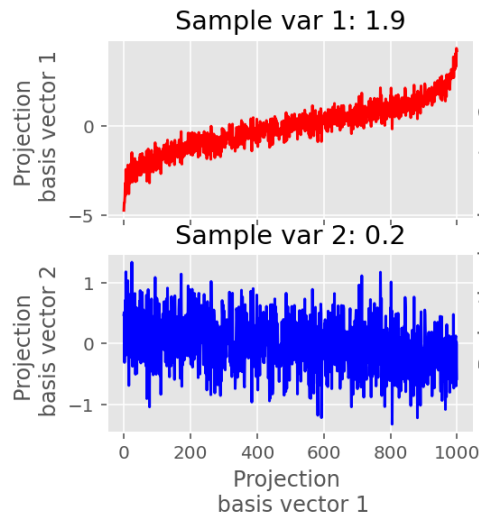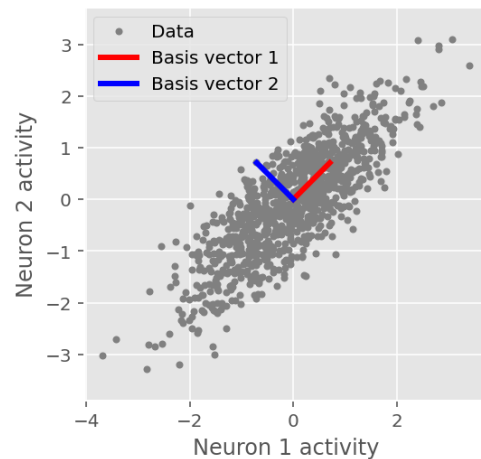Samples   **X**   ➡️   Samples   **?**

# Covariance reveals structure



**Expectation:**
- Directions of large variation represent signal.
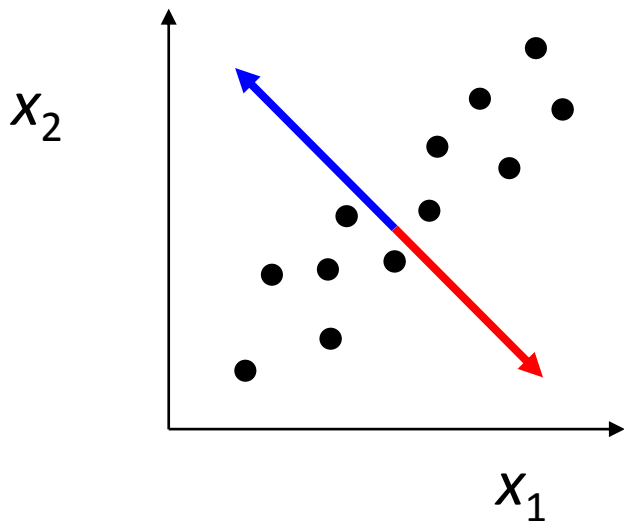- Directions of small variation represent noise.

**Goal of PCA:** Find directions of maximum variance.

# Projected variance is largest when basis is aligned with covariance direction

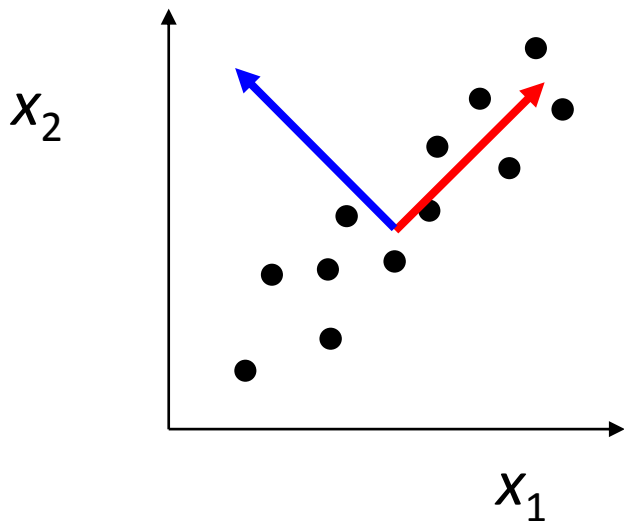# Covariance reveals structure



**Expectation:**
- Directions of large variation represent signal.
- Directions of small variation represent noise.

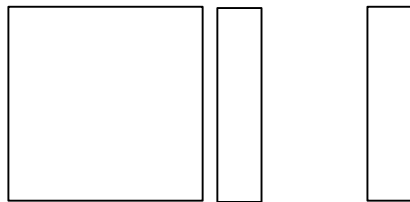**Goal of PCA:** Find directions of maximum variance.
- $\mathbf{w}_1$ – vector that has highest projected variance
- $\mathbf{w}_2$ – vector that is orthogonal to $\mathbf{w}_1$ and has highest projected variance
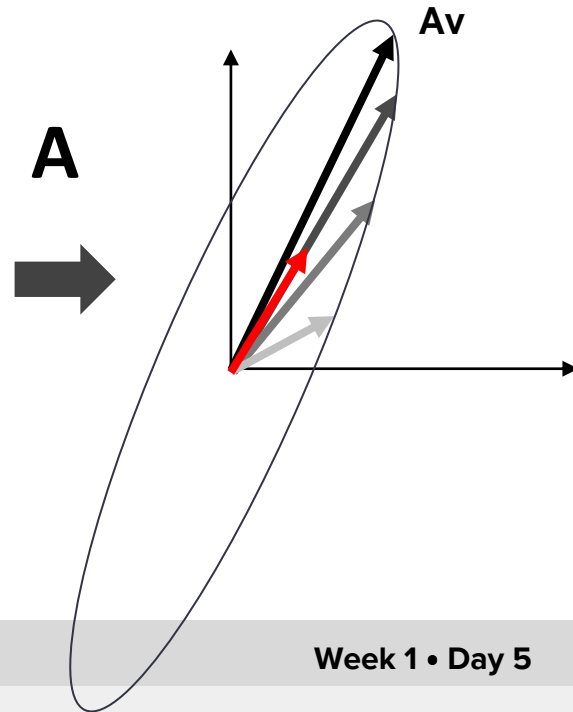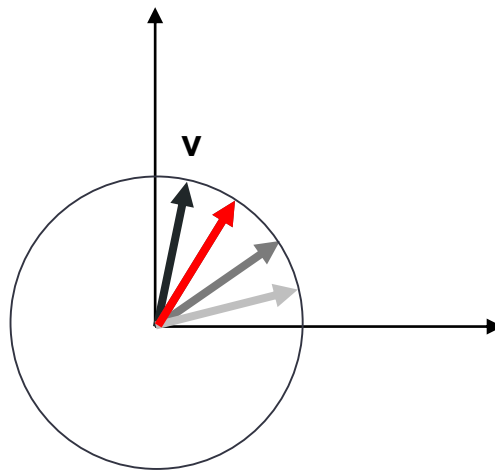- Etc.

# Covariance reveals structure



**Expectation:**
- Directions of large variation represent signal.
- Directions of small variation represent noise.

**Goal of PCA:** Find directions of maximum variance.
- $w_1$ – vector that has highest projected variance
- $w_2$ – vector that is orthogonal to $w_1$ and has highest projected variance
- Etc.

**Mathematical solution:**
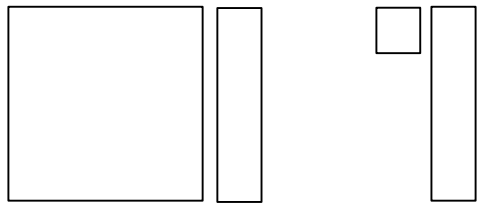- $w_1$, $w_2$, …, $w_N$ are the eigenvectors of $\Sigma$.

# Eigenvalue refresher

A matrix is a linear transformation

**A    v   =   ?**

# Eigenvalue refresher

A matrix is a linear transformation

$$\mathbf{A} \quad \textcolor{red}{\mathbf{v}} \; = \; \textcolor{blue}{\lambda} \; \textcolor{red}{\mathbf{v}}$$

$\textcolor{blue}{\lambda}$: eigenvalue
$\textcolor{red}{\mathbf{v}}$: eigenvector

**A**

# Covariance reveals structure



**Expectation:**
- Directions of large variation represent signal.
- Directions of small variation represent noise.

**Goal of PCA:** Find directions of maximum variance.
- $\mathbf{w}_1$ – vector that has highest projected variance
- $\mathbf{w}_2$ – vector that is orthogonal to $\mathbf{w}_1$ and has highest projected variance
- Etc.

**Mathematical solution:**
- $\mathbf{w}_1, \mathbf{w}_2, ..., \mathbf{w}_N$ are the eigenvectors of $\mathbf{\Sigma}$.
- Projected variance onto each $\mathbf{w}_i$ is given by its corresponding eigenvalue $\lambda_i$.
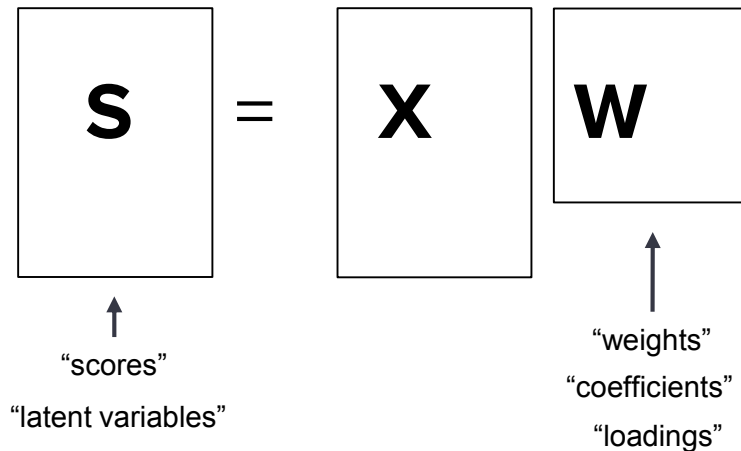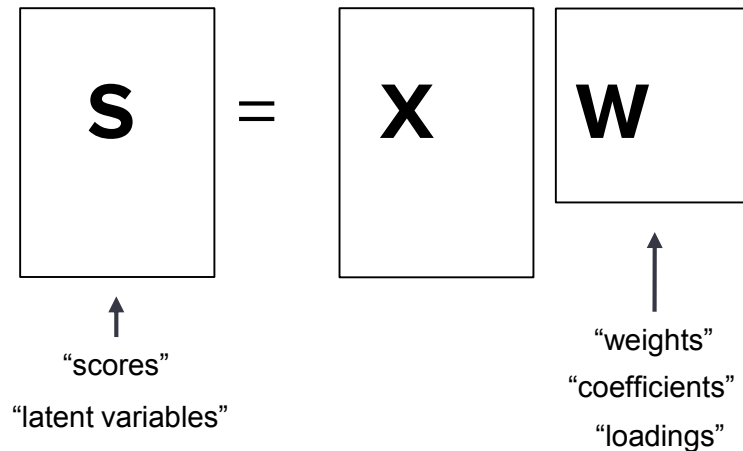
# How to perform PCA

**Basic algorithm**
1. Subtract the mean
2. Calculate the eigenvectors $\mathbf{w}_i$ of the covariance matrix $\Sigma$, ordered by their corresponding eigenvalue $\lambda_i$.
3. Project the data $\mathbf{X}$ onto the new basis $\mathbf{w}_1, \mathbf{w}_2, ..., \mathbf{w}_N$.

**Key properties**
- $\mathbf{w}_i$ are orthogonal
- $\mathbf{s}_i$ are uncorrelated
- Projected variance = $\lambda_i$

$$\mathbf{S} = \mathbf{X} \cdot \mathbf{W}$$

"scores"
"latent variables"

"weights"
"coefficients"
"loadings"

# (break for tutorial exercise)

neuromatch
academy

# How to perform PCA

**Basic algorithm**
1. Subtract the mean
2. Calculate the eigenvectors $\mathbf{w}_i$ of the covariance matrix $\mathbf{\Sigma}$, ordered by their corresponding eigenvalue $\lambda_i$.
3. Project the data $\mathbf{X}$ onto the new basis $\mathbf{w}_1, \mathbf{w}_2, ..., \mathbf{w}_N$.

**Key properties**
- $\mathbf{w}_i$ are orthogonal
- $\mathbf{s}_i$ are uncorrelated
- Projected variance = $\lambda_i$

$$\boxed{\mathbf{S}} = \boxed{\mathbf{X}}\ \boxed{\mathbf{W}}$$

"scores"

"latent variables"

"weights"

"coefficients"

"loadings"

# $\mathbf{w}_i$ are orthogonal

**What we know**
- Since $\mathbf{w}_i$ eigenvectors of the covariance matrix

$$\hat{\boldsymbol{\Sigma}}\mathbf{w}_i = \lambda_i \mathbf{w}_i$$

- Covariance matrix is symmetric

$$\hat{\boldsymbol{\Sigma}} = \hat{\boldsymbol{\Sigma}}^T$$

$$\lambda_j \mathbf{w}_i^T \mathbf{w}_j$$

**What we want to show**

$$\mathbf{w}_i \cdot \mathbf{w}_j = 0$$

$$(\lambda_j - \lambda_i)\mathbf{w}_i \cdot \mathbf{w}_j = 0$$

**If the eigenvalues are different, then the eigenvectors must be orthogonal.**

# $\mathbf{s}_i$ are uncorrelated

**What we know**

- Since $\mathbf{w}_i$ eigenvectors of the covariance matrix
$$\hat{\boldsymbol{\Sigma}}\mathbf{w}_i = \lambda_i \mathbf{w}_i$$

- Scores $\mathbf{s}_i$ represent the projected data
$$\mathbf{s}_i = \mathbf{X}\mathbf{w}_i$$

- Since $\mathbf{X}$ is zero-mean, $\mathbf{s}_i$ is zero-mean
$$\bar{\mathbf{s}}_i = 0$$

$$\mathrm{cov}(\mathbf{s}_i, \mathbf{s}_j)$$

**What we want to show**
$$\mathrm{cov}(\mathbf{s}_i, \mathbf{s}_j) = 0$$

The scores (projected data) are uncorrelated.

# $\lambda_i$ describe projected variances

**What we know**
- Since $\mathbf{w}_i$ eigenvectors of the covariance matrix
$$\hat{\boldsymbol{\Sigma}}\mathbf{w}_i = \lambda_i\mathbf{w}_i$$
- Scores $\mathbf{s}_i$ represent the projected data
$$\mathbf{s}_i = \mathbf{X}\mathbf{w}_i$$
- Since **X** is zero-mean, $\mathbf{s}_i$ is zero-mean
$$\bar{\mathbf{s}}_i = 0$$

**What we want to show**
$$\mathrm{var}(\mathbf{s}_i) = \lambda_i$$

$$\mathrm{var}(\mathbf{s}_i) = \frac{1}{N_{\mathrm{samples}}}\mathbf{s}_i^T\mathbf{s}_i - \bar{\mathbf{s}}_i^2$$

$$= \frac{1}{N_{\mathrm{samples}}}\mathbf{s}_i^T\mathbf{s}_i$$

$$= \frac{1}{N_{\mathrm{samples}}}(\mathbf{X}\mathbf{w}_i)^T\mathbf{X}\mathbf{w}_i$$

$$= \frac{1}{N_{\mathrm{samples}}}\mathbf{w}_i^T\mathbf{X}^T\mathbf{X}\mathbf{w}_i$$

$$= \mathbf{w}_i^T\hat{\boldsymbol{\Sigma}}\mathbf{w}_i$$

$$= \lambda_i\mathbf{w}_i^T\mathbf{w}_i$$

**The variance of the scores (projected data) is its corresponding eigenvalue.**
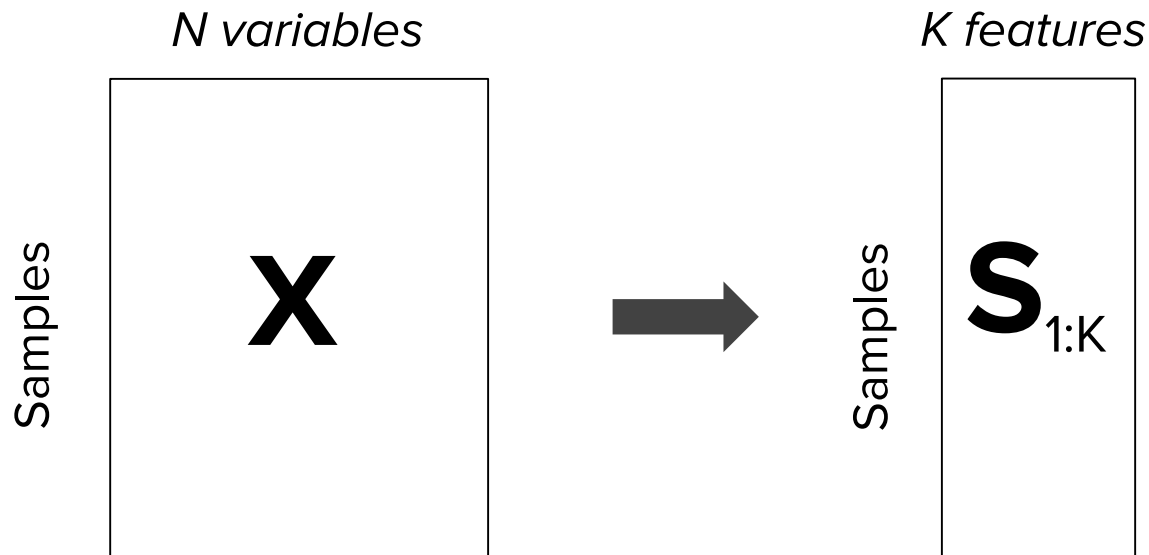
# Dimensionality reduction and reconstruction

Tutorial 3

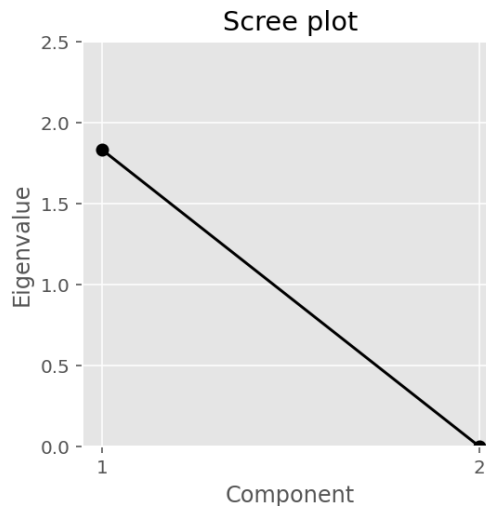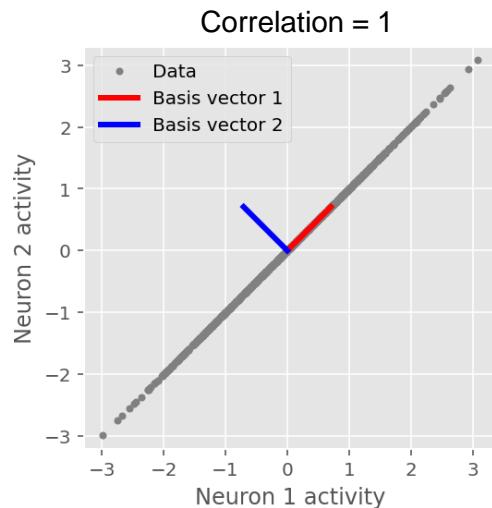# Goal of dimensionality reduction

*N variables*

*K features*

Samples

**X**

→

Samples

**?**

# Dimensionality reduction via PCA



*N variables*

Samples

**X**

*K features*

Samples

**S**$_{1:K}$

# Intrinsic vs. extrinsic dimensionality

**Correlation = 1**

Data
Basis vector 1
Basis vector 2

Neuron 2 activity

Neuron 1 activity

**Scree plot**

Eigenvalue

Component

**Extrinsic** dimensionality:
N = 2

**Intrinsic** dimensionality:
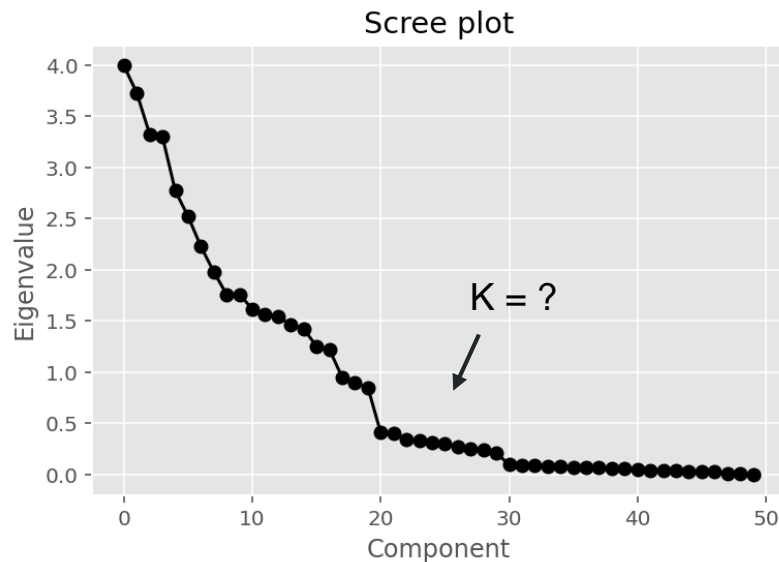K = 1

# Intrinsic vs. extrinsic dimensionality



**Extrinsic**
dimensionality:
N = 2

**Intrinsic**
dimensionality:
K = 2

# Intrinsic vs. extrinsic dimensionality



**Extrinsic** dimensionality:
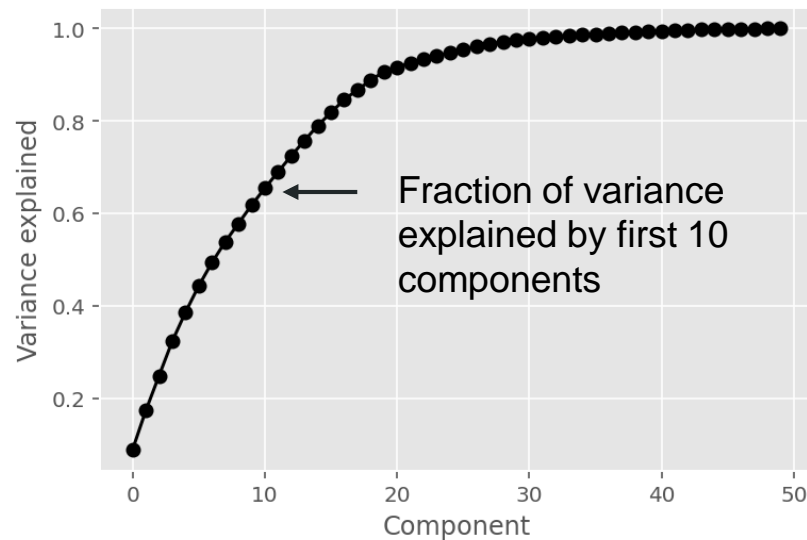
N = 2

**Intrinsic** dimensionality:

K = ?

# How to determine intrinsic dimensionality?
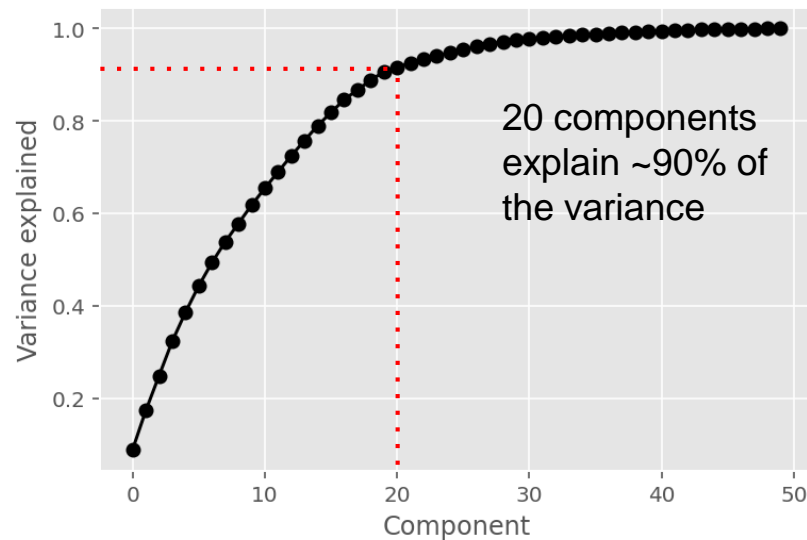


Scree plot

K = 20
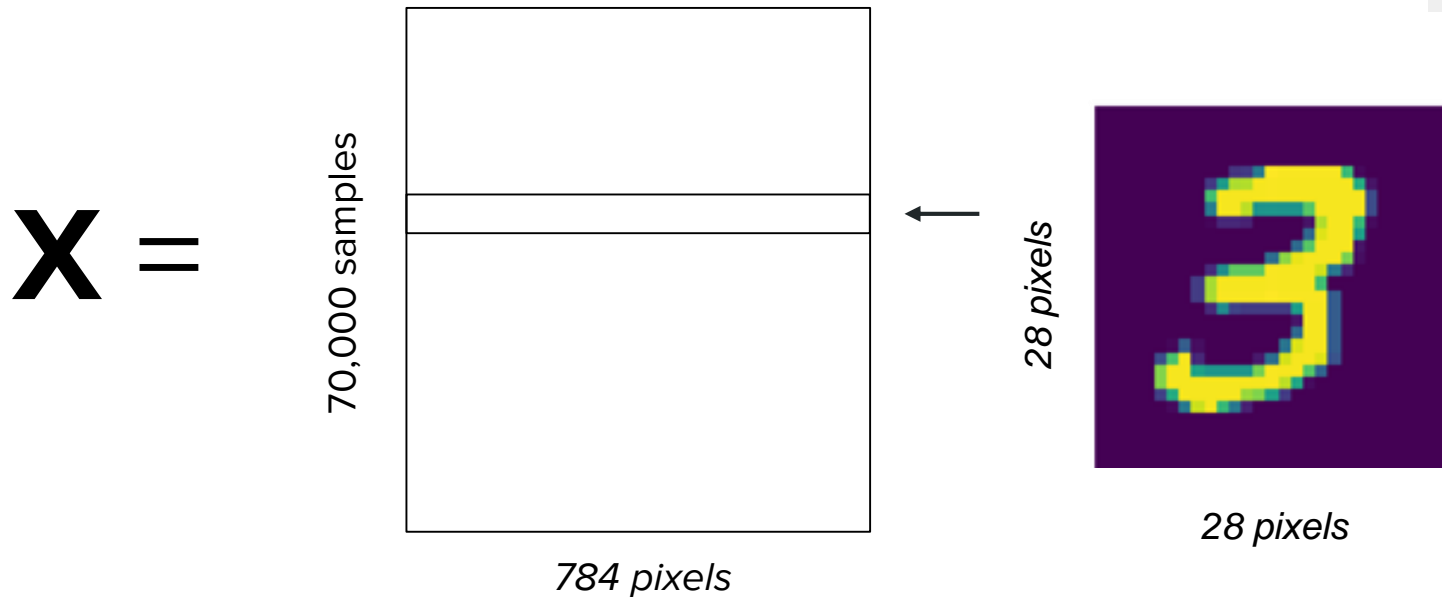
# How to determine intrinsic dimensionality?



Scree plot

K = ?

# Total variance explained



Fraction of variance explained by first 10 components

# Total variance explained



20 components explain ~90% of the variance

# The MNIST dataset

# The MNIST dataset

$$\mathbf{X} =$$



70,000 samples

784 pixels

28 pixels

28 pixels

# Reconstruction from PCA

Once we have **S** and **W** how do we reconstruct **X** ?

$$\boxed{\textbf{S}} = \boxed{\textbf{X}} \boxed{\textbf{W}}$$

**Algorithm for PCA**
1. Subtract the mean
2. Calculate the eigenvectors $\mathbf{w}_i$ of the covariance matrix $\mathbf{\Sigma}$, ordered by their corresponding eigenvalue $\lambda_i$.
3. Project the data **X** onto the new basis $\mathbf{w}_1, \mathbf{w}_2, ..., \mathbf{w}_N$.

# Reconstruction from PCA

Once we have **S** and **W** how do we reconstruct **X** ?

$$\boxed{\textbf{S}} \; \boxed{\textbf{W}^{\text{T}}} = \boxed{\textbf{X}} \; \underbrace{\boxed{\textbf{W}} \; \boxed{\textbf{W}^{\text{T}}}}$$

= Identity matrix

because $\textbf{w}_i$ are
orthonormal basis

# Reconstruction from PCA

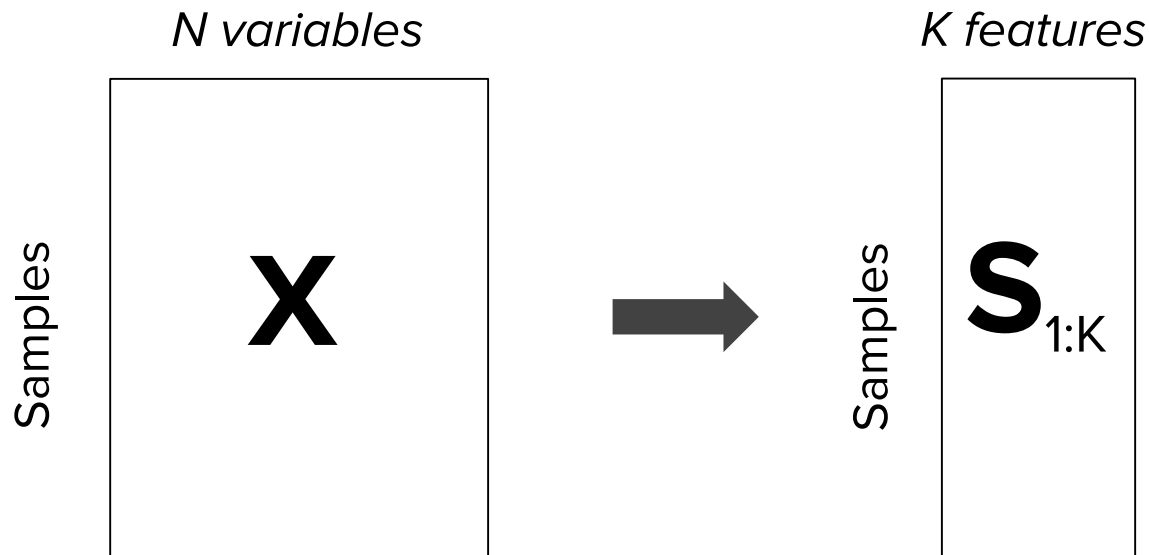Once we have **S** and **W** how do we reconstruct **X** ?

$$ \mathbf{S} \quad \mathbf{W}^{\mathrm{T}} = \mathbf{X} $$

**Algorithm for reconstruction from PCA**
1. Multiply scores by transpose of the weight matrix
2. Add the mean

# Dimensionality reduction via PCA

*N variables*

*K features*

Samples

**X**

Samples

**S**$_{1:K}$

# Reconstruction from PCA

Once we have **S** and **W** how do we reconstruct **X** ?

$$\mathbf{S}_{1:K} \ (\mathbf{W}_{1:K})^{\mathsf{T}} = \hat{\mathbf{X}}$$

**Algorithm for reconstruction from PCA**
1. Truncate scores and weight matrix after top K components
2. Multiply scores by transpose of the weight matrix
3. Add the mean

**Goal of PCA:** Find K-dimensional basis that minimizes the reconstruction error.

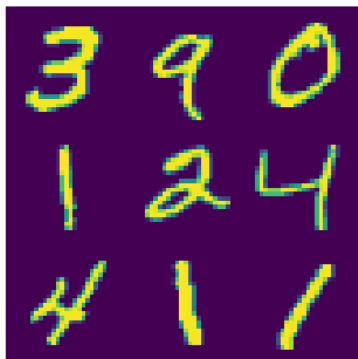$$\|\mathbf{X} - \hat{\mathbf{X}}\|^2$$

# Nonlinear dimensionality reduction

### Tutorial 4

# PCA: the big picture



Latent subspace

$x_2$
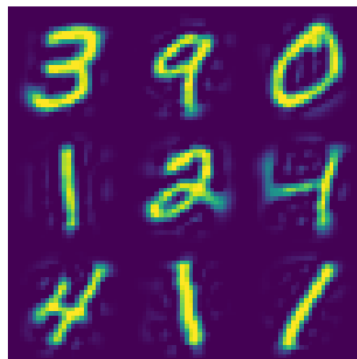
$x_3$
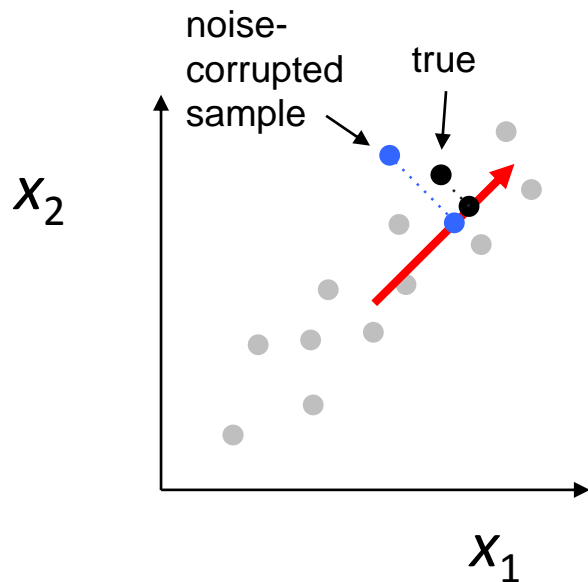
$x_1$

Latent representation

Component 2

Component 1

# PCA for compression

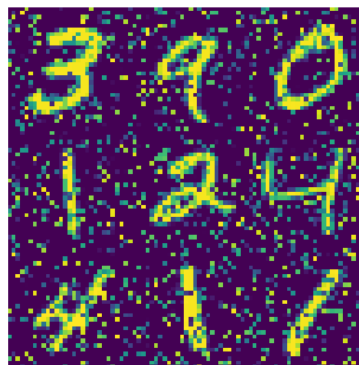Data

Reconstructed, K=71

$$\mathbf{S}_{1:K} \; (\mathbf{W}_{1:K})^{\mathsf{T}} \approx \mathbf{X}$$
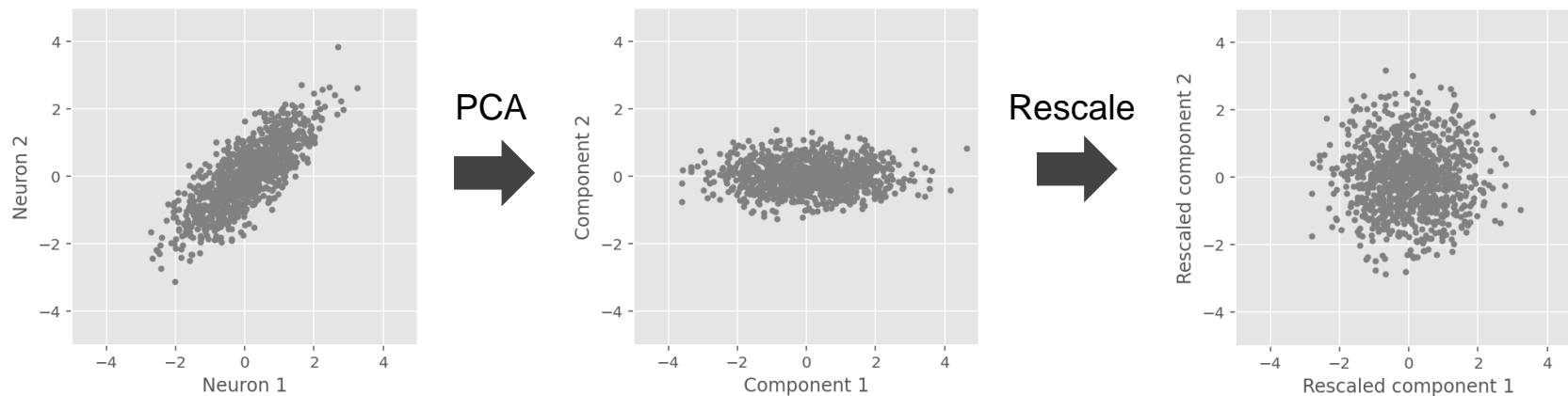
# PCA for denoising



noise-corrupted sample

true

$x_2$

$x_1$

Data
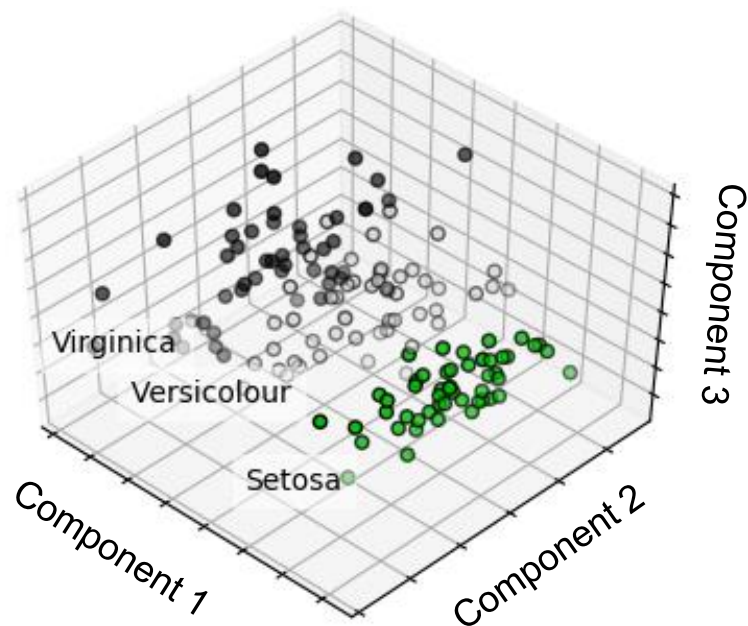
Reconstructed

# PCA for whitening



PCA

Rescale

# PCA for visualization



Source: scikit-learn.org

# Visualizing MNIST

# (break for tutorial exercise)

neuromatch
academy

# Linear dimensionality reduction

Linear transformation to lower
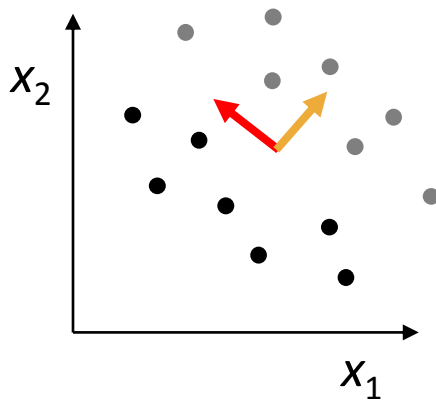dimensional representation **Y**

$$\boxed{\mathbf{Y}} = \boxed{\mathbf{X}} \boxed{\mathbf{W}}$$

**Probabilistic PCA (PPCA):**
- Explicit noise model $\quad \mathbf{x}|\mathbf{y} \sim \mathcal{N}(\mathbf{yW}, \sigma_\epsilon^2 \mathbf{I})$
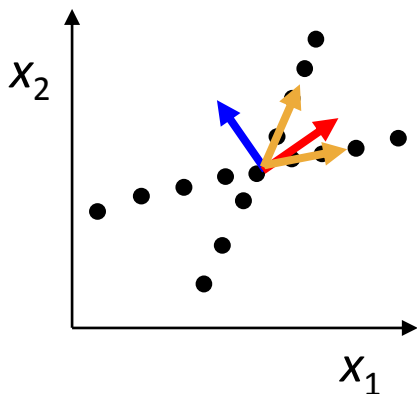
**Factor Analysis (FA):**
- Non-isotropic noise $\quad \mathbf{x}|\mathbf{y} \sim \mathcal{N}(\mathbf{yW}, \mathbf{D})$
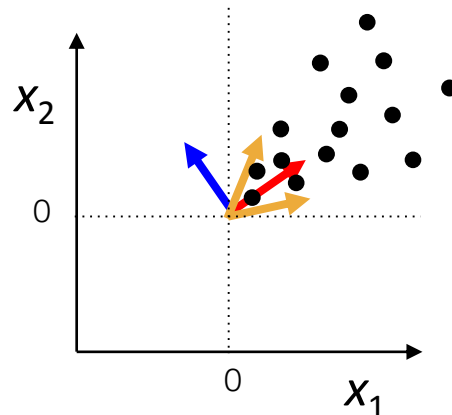
$x_2$

$x_1$

**Linear discriminant analysis (LDA) :**
- Preserve class discriminatory information
- Example of **supervised** dimensionality reduction

# Blind source separation
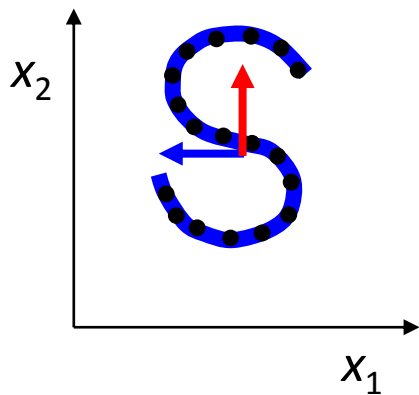


**Independent Components Analysis (ICA) :**
- Stronger condition than uncorrelated
- Basis vectors not necessarily orthogonal
- Components not ordered by importance

**Nonnegative Matrix Factorization (NMF) :**
- Weights and components positive
- Basis vectors not necessarily orthogonal
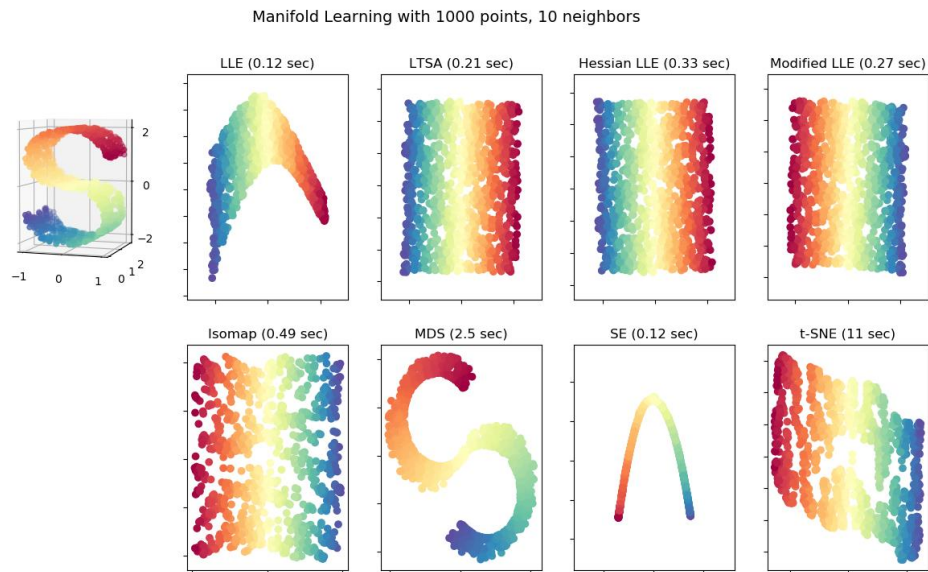- No linear mapping to low-D space

# When is linear not enough?

$x_2$

$x_1$

"embedding"

# Nonlinear dimensionality reduction



Manifold Learning with 1000 points, 10 neighbors
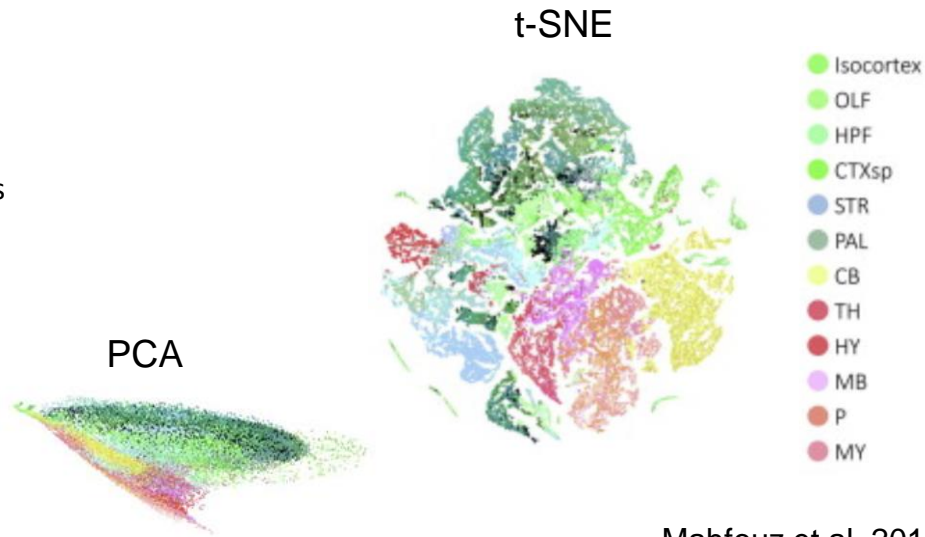
Source: scikit-learn.org

# t-distributed stochastic neighbor embedding

**t-SNE**
- Visualization in 2D or 3D
- Define similarity between samples **X**
- Find a mapping to low-dimensional **Y** that preserves similarities as much as possible

**Differences from PCA**
- Nonlinear
- Stochastic
- No reconstruction
- Free parameter: *perplexity*

t-SNE

PCA

- Isocortex
- OLF
- HPF
- CTXsp
- STR
- PAL
- CB
- TH
- HY
- MB
- P
- MY

Mahfouz et al. 2015

# (break for tutorial exercise)

neuromatch
academy