

# Lecture 14

# Principal Component Analysis

Jochen Braun

Otto-von-Guericke-Universität Magdeburg,  
Cognitive Biology Group

Engineering Neuroscience / Computational Neuroscience II  
SS 2020

Credits: Jonathan Shlens (2005) “A Tutorial on Principal  
Component Analysis”

## 14. Principal Component Analysis (PCA)

*PCA is a useful way to summarize high-dimensional data (repeated observations of multiple variables). This lecture provides the underlying linear algebra needed for practical applications. It also emphasizes consistent notation. The central ideas of PCA are **orthonormal coordinate** systems, the distinction between **variance and covariance**, and the possibility of choosing an orthonormal basis to **eliminate covariance**. Technically, PCA may be performed either by **eigenvector analysis** of the covariance matrix or by **singular value decomposition** of the original observation matrix. Both variants are described. We illustrate the method with a non-biological example (image-denoising) and with a biological example (multi-unit activity, see **Exercise 8**).*

# Overview

- ▶ **1 British food**
- ▶ **2 Orthonormal bases**
- ▶ **3 Variance and covariance**
- ▶ **4 Diagonalizing variance**
- ▶ **5 De-noising**
- ▶ **6 Multi-unit activity**
  
- ▶ **Appendix: Singular value decomposition**

# Motivation

Principal component analysis (PCA) is one of the most valuable results of applied linear algebra.

It is widely used – from neuroscience to computer graphics – because it is an easy way to simplify confusing data sets. With minimal effort, PCA dramatically reduces the dimensionality of a large data set and potentially reveals a simpler structure hiding within.

During brain development, an (approximate) PCA of sensory inputs is performed by activity-dependent Hebbian plasticity.

# Variables and observations

We distinguish between *variables* (types of measurements) and *observations* (times or conditions of measurements).

Think of variables as 'effects' (possibly mixed, redundant, or noisy) and observations as variations over 'causes' (due to differences in time, condition, or other factors).

For example, what kind of food people eat is a *variable*, in that it is an 'effect' of many causal factors such as climate, culture, economics, and so on.

When considering different countries, historical periods, or climate zones, these causal factors will presumably differ from case to case. The different cases being compared may be considered *observations*.

# 1. British food consumption

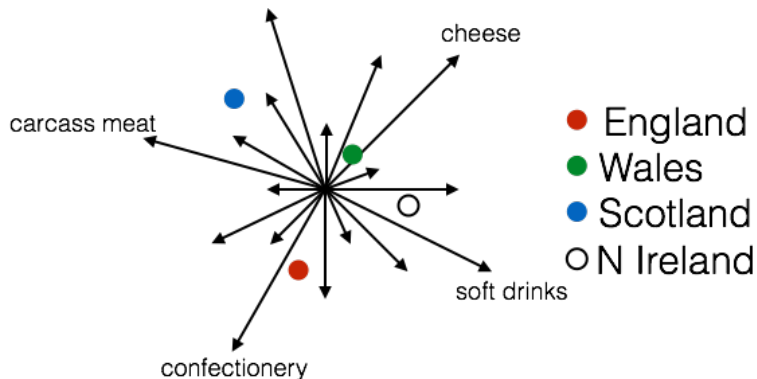
Consumption of different foods in Britain, in g/person/week:

	Eng	Wal	Scot	N Ire
cheese	105	103	103	66
carcass meat	245	227	242	267
other meat	685	803	750	586
fish	147	160	122	93
fats and oils	193	235	184	209
sugars	156	175	147	139
fresh potatoes	720	874	566	1033
fresh veg	253	265	171	143
other veg	488	570	418	355
processed potatoes	198	203	220	187
processed veg	360	365	337	334
fresh fruit	1102	1137	957	674
cereals	1472	1582	1462	1494
hot beverages	57	73	53	47
soft drinks	1374	1256	1572	1506
alcoholic drinks	375	475	458	135
confectionery	54	64	62	41

We consider foods as *variables* and countries as *observations*!  
(Credits: Mark Richardson).

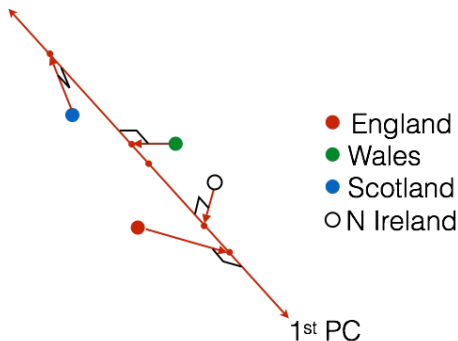
## 4 points in 17-dimensional space

We form an orthonormal basis of 17 food types and represent 4 observations of 17 variables as four points in this space (one point per observation):



## Project to special line

We may project our observations orthogonally on any line in 17D space. We seek the particular line along which the projections scatter most widely (show maximal variance).

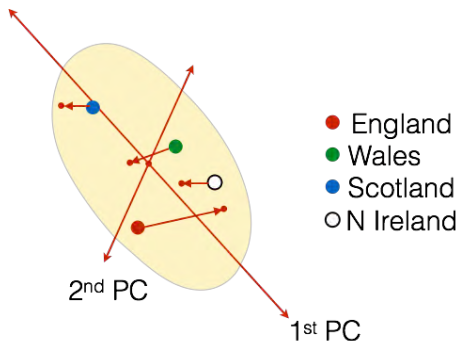


We term this special line the *1<sup>st</sup> principal component* .



## Project to special plane

Similarly, we may project orthogonally on any plane in 17D space. Again we seek a particular plane: the plane that spans the two orthogonal lines along which the projections scatter most widely and second most widely.



We term these orthogonal lines the *1<sup>st</sup> and 2<sup>nd</sup> principal component* .

# Maximal variance between observations (countries)

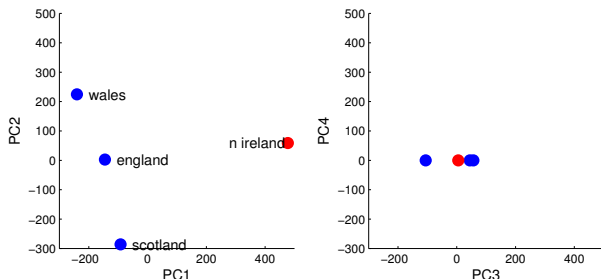
The axis of maximal variance (1st principal component) between countries is here shown horizontally



From this perspective, it is evident that England, Wales, and Scotland (blue) are comparatively similar, whereas N Ireland (red) is comparatively distinct.

# Principal components compared

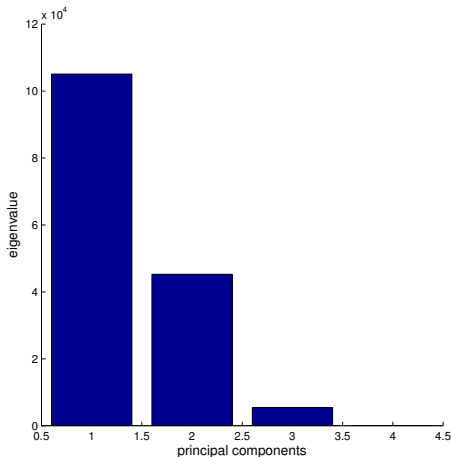
In total, there are 17 'principal components'. Together, they form an alternative orthonormal basis for our space. The 1<sup>st</sup> and 2<sup>nd</sup> principal components are shown on the left, the 3<sup>rd</sup> and 4<sup>th</sup> on the right:



Note that the first two PCs capture most of the variance between observations (countries).

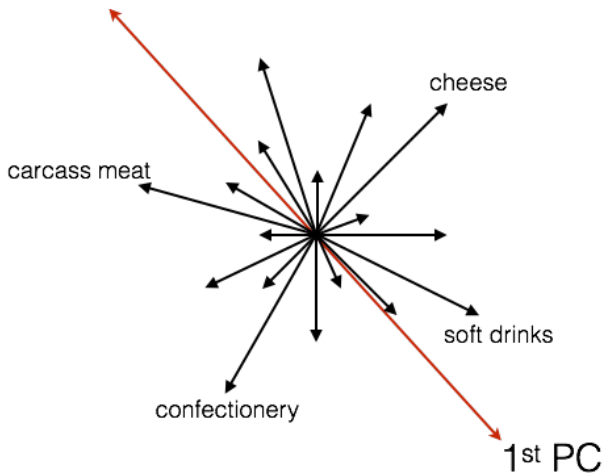
## Variance captured

The relative importance of different principal components is measured by the variance captured. In our example, only four principal components capture variance. (This is because there were only 4 observations).



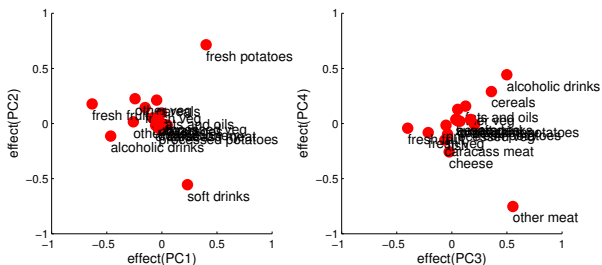
## Contributions to principal component

A 'principal component' is linear combination of variables. In our example, it is a direction in the 17D space of food types.



# Visualizing effects

To visualize the contributions ('effects') of different variables (foods) to a principal component, we can project the original 17 orthonormal basis vectors onto the PC of interest:



The largest effect on PC1 have fresh fruit and alcoholic drinks (both of which are consumed less in N Ireland). The largest effect on PC2 have fresh potatoes and soft drinks.

# Summary

- ▶ In many fields (including neuroscience), we observe a multitude of variables under different conditions or at different times.
- ▶ It often makes sense to seek a simpler way to describe of the total variance (sum of variances over condition or time).
- ▶ PCA considers observations (conditions or times) as points in a high-dimensional space, the dimensionality of which is determined by the number of variables.
- ▶ PCA computes an alternative orthonormal basis for this space and it ranks the basis vectors ('principal components') by the variance captured.
- ▶ In effect, it identifies the points of view (directions of projection) from which observations appear most variable.

[www.models.life.ku.dk](http://www.models.life.ku.dk)

<https://www.youtube.com/user/QualityAndTechnology/playlists>  
Principal Component Analysis – PCA

<https://www.youtube.com/watch?v=K-F19DORO1w&list=PLBC24FD8C389FE9E4>

8:20 to 10:20



## 2 Orthonormal bases

Consider  $m$  variables, **each with zero mean** and observed  $n$  times, with  $i^{\text{th}}$  observations grouped into **column** vectors (for now!)

$$\hat{\mathbf{x}}_i = \begin{pmatrix} x_{1i} \\ \vdots \\ x_{mi} \end{pmatrix} \quad \sum_{i=1}^n x_{ji} = 0, \quad i = 1 \dots n, \quad j = 1 \dots m$$

and all  $n$  observations collected into a matrix with  $N$  columns and  $M$  rows

$$\mathbf{X} = (\hat{\mathbf{x}}_1 \quad \hat{\mathbf{x}}_2 \quad \dots \quad \hat{\mathbf{x}}_n) = \begin{pmatrix} x_{11} & \dots & x_{1n} \\ \vdots & \ddots & \vdots \\ x_{m1} & \dots & x_{mn} \end{pmatrix} \in \mathbb{R}^{m \times n}$$

Note the indexing with  $x_{\text{row}, \text{column}}$  Or  $x_{\text{variable}, \text{observation}}$ :

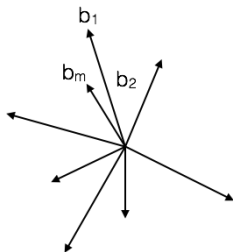
**rows**  $\simeq$  **variables** and **columns**  $\simeq$  **observations**

## Orthonormal basis

Implicitly, this assumes a basis of  $m$  orthonormal **row vectors**

$\mathbf{b}_j = (b_{11} \dots b_{1m})$ , each with  $m$  elements,  $\mathbf{B} \in \mathbb{R}^{m \times m}$ :

$$\mathbf{B} = \begin{pmatrix} \mathbf{b}_1 \\ \mathbf{b}_2 \\ \vdots \\ \mathbf{b}_m \end{pmatrix} = \begin{pmatrix} b_{11} & b_{12} & \dots & b_{1m} \\ b_{21} & b_{22} & \dots & b_{2m} \\ \vdots & \vdots & \ddots & \vdots \\ b_{m1} & b_{m2} & \dots & b_{mm} \end{pmatrix} = \begin{pmatrix} 1 & 0 & \dots & 0 \\ 0 & 1 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & 1 \end{pmatrix} = \mathbf{I}$$



'Orthonormal' means of unit length and pairwise orthogonal

$$\mathbf{b}_j \cdot \mathbf{b}_k = \delta_{jk}$$

$$\mathbf{B}\mathbf{B}^T = \mathbf{B}^T\mathbf{B} = \mathbf{I}$$

Project **columns**  $\hat{\mathbf{x}}_i \in \mathbb{R}^{m \times 1}$  of *individual* observations onto **rows**  $\mathbf{b}_j$  of basis

$$\hat{\mathbf{x}}_i = \mathbf{B}\hat{\mathbf{x}}_i = \begin{pmatrix} \mathbf{b}_1 \cdot \hat{\mathbf{x}}_i \\ \mathbf{b}_2 \cdot \hat{\mathbf{x}}_i \\ \vdots \\ \mathbf{b}_m \cdot \hat{\mathbf{x}}_i \end{pmatrix} = \begin{pmatrix} b_{11}x_{1i} + \dots + b_{1m}x_{mi} \\ b_{21}x_{1i} + \dots + b_{2m}x_{mi} \\ \vdots \\ b_{m1}x_{1i} + \dots + b_{mm}x_{mi} \end{pmatrix} = \begin{pmatrix} x_{1i} \\ x_{2i} \\ \vdots \\ x_{mi} \end{pmatrix}$$

Project entire matrix  $\mathbf{X} \in \mathbb{R}^{m \times n}$  of *all* observations onto basis

$$\mathbf{X} = \mathbf{B}\mathbf{X} = (\mathbf{B}\hat{\mathbf{x}}_1 \dots \mathbf{B}\hat{\mathbf{x}}_n) = \begin{pmatrix} \mathbf{b}_1 \cdot \hat{\mathbf{x}}_1 & \dots & \mathbf{b}_1 \cdot \hat{\mathbf{x}}_n \\ \mathbf{b}_2 \cdot \hat{\mathbf{x}}_1 & \dots & \mathbf{b}_2 \cdot \hat{\mathbf{x}}_n \\ \vdots & \ddots & \dots \\ \mathbf{b}_m \cdot \hat{\mathbf{x}}_1 & \dots & \mathbf{b}_m \cdot \hat{\mathbf{x}}_n \end{pmatrix} \in \mathbb{R}^{m \times n}$$

# Many orthonormal bases are possible

- ▶ The orthonormal basis  $\mathbf{B} \in \mathbb{R}^{m \times m}$  reflects how we happened to collect the observations.
- ▶ Many other orthonormal bases  $\mathbf{P} \in \mathbb{R}^{m \times m}$  may be obtained as linear combinations of  $\mathbf{B}$ .
- ▶ Perhaps in terms of other orthonormal bases, we can re-express our observations more simply or more meaningfully?

## Alternative basis

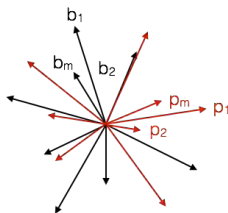
Alternative basis  $\mathbf{P} \in \mathbb{R}^{m \times m}$  of orthonormal **row** vectors

$\mathbf{p}_j = (p_{j1} \dots p_{jm})$ :

$$\mathbf{P} = \begin{pmatrix} \mathbf{p}_1 \\ \vdots \\ \mathbf{p}_m \end{pmatrix} = \begin{pmatrix} p_{11} & \dots & p_{1m} \\ \vdots & \ddots & \vdots \\ p_{m1} & \dots & p_{mm} \end{pmatrix} \in \mathbb{R}^{m \times m}$$

$$\mathbf{p}_j \cdot \mathbf{p}_k = \delta_{jk}$$

$$\mathbf{P}\mathbf{P}^T = \mathbf{P}^T\mathbf{P} = \mathbf{I}$$



## Projection to alternative basis

Transform individual observations  $\mathbf{x}_i \rightarrow \mathbf{y}_i!$  Obtain  $y$ 's as dot products (projections) of  $\mathbf{x}_i$  on  $\mathbf{p}_j!$

$$\hat{\mathbf{y}}_i = \mathbf{P}\hat{\mathbf{x}}_i = \begin{pmatrix} \mathbf{p}_1 \cdot \hat{\mathbf{x}}_i \\ \mathbf{p}_2 \cdot \hat{\mathbf{x}}_i \\ \vdots \\ \mathbf{p}_m \cdot \hat{\mathbf{x}}_i \end{pmatrix} = \begin{pmatrix} p_{11}x_{1i} + \dots + p_{1m}x_{mi} \\ p_{21}x_{1i} + \dots + p_{2m}x_{mi} \\ \vdots \\ p_{m1}x_{1i} + \dots + p_{mm}x_{mi} \end{pmatrix} = \begin{pmatrix} y_{1i} \\ y_{2i} \\ \vdots \\ y_{mi} \end{pmatrix}$$

Transform all observations  $\mathbf{X} \rightarrow \mathbf{Y}$

$$\mathbf{Y} = \mathbf{P}\mathbf{X} = (\mathbf{P}\hat{\mathbf{x}}_1 \dots \mathbf{P}\hat{\mathbf{x}}_m) = \begin{pmatrix} \mathbf{p}_1 \cdot \hat{\mathbf{x}}_1 & \dots & \mathbf{p}_1 \cdot \hat{\mathbf{x}}_m \\ \mathbf{p}_2 \cdot \hat{\mathbf{x}}_1 & \dots & \mathbf{p}_2 \cdot \hat{\mathbf{x}}_m \\ \vdots & \ddots & \dots \\ \mathbf{p}_m \cdot \hat{\mathbf{x}}_1 & \dots & \mathbf{p}_m \cdot \hat{\mathbf{x}}_m \end{pmatrix} \in \mathbb{R}^{m \times n}$$

Keep in mind that  $\hat{\mathbf{x}}_i$  and  $\hat{\mathbf{y}}_i$  are columns,  $\mathbf{b}_i$  and  $\mathbf{p}_i$  are rows

# Summary

- ▶ Observations  $\mathbf{X}$  need not be expressed in native basis  $\mathbf{B}$ .
- ▶ Alternative orthonormal bases  $\mathbf{P}$  are linear combinations of  $\mathbf{B}$ .
- ▶ Transformation to new basis is performed conveniently by matrix multiplication

$$\mathbf{Y} = \mathbf{P} \mathbf{X}$$

because dot products compute projections!

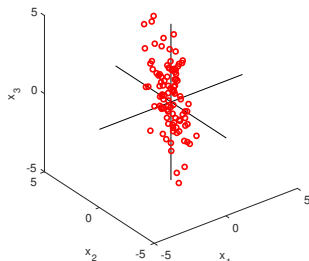
- ▶ Row vectors  $\mathbf{p}_j$  will become *principal components* of  $\mathbf{X}$ .
- ▶ What would be a good choice for  $\mathbf{P}$ ?
- ▶ In which respect could  $\mathbf{Y}$  be an improvement over  $\mathbf{X}$ ?

### 3 Variance

Consider the squared norm of individual observations  $\hat{\mathbf{x}}_i$  (= square distance from origin) and the sum over all observations  $\mathbf{X}$

$$|\hat{\mathbf{x}}_i|^2 = x_{1i}^2 + \dots + x_{mi}^2 = \sum_{j=1}^m x_{ji}^2, \quad \sum_i |\hat{\mathbf{x}}_i|^2 = \sum_{i=1}^n \sum_{j=1}^m x_{ji}^2$$

Observations in 3D



Neither the squared norm of individual observations nor their sum over all observations depend on the chosen basis!

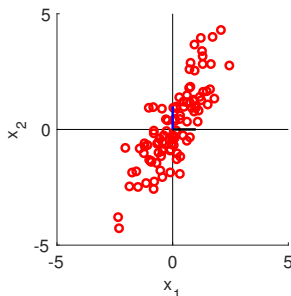


## Variance within rows

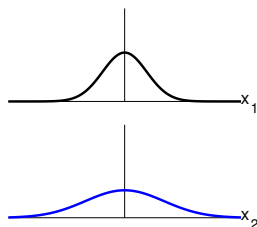
Define  $\mathbf{x}_j$  (distinct from  $\hat{\mathbf{x}}_j$ !) to represent **zero-mean rows** of  $x_{ji}$

$$\mathbf{X} = \begin{pmatrix} x_{11} & \dots & x_{1n} \\ \vdots & \ddots & \vdots \\ x_{m1} & \dots & x_{mn} \end{pmatrix} = \begin{pmatrix} \mathbf{x}_1 \\ \vdots \\ \mathbf{x}_m \end{pmatrix} \quad \begin{array}{l} \sum_i x_{1i} = 0 \\ \vdots \\ \sum_i x_{mi} = 0 \end{array}$$

Typically, there is some, but not necessarily the same, variance within each of the different rows or dimensions, here  $x_1$  and  $x_2$ :

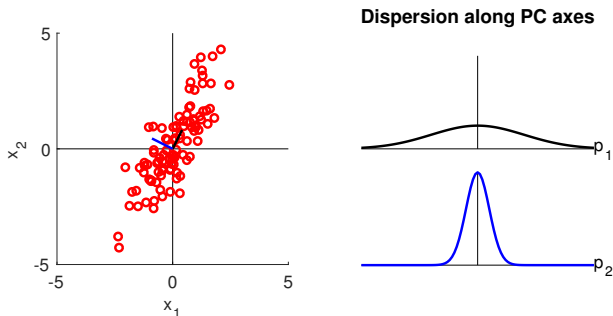


Dispersion along original axes



# Variance and orthonormal basis

The variance within each row or dimension typically changes with the basis. Importantly, the total variance remains always the same!



In the illustration, the rotated basis *maximizes* variance in one dimension ( $p_1$ ) and minimizes it in the other ( $p_2$ ).

## Variance definition

The variance of **zero-mean** variable  $j$  over observations  $i = 1, \dots, n$  is defined as

$$\begin{aligned}\sigma_j^2 &= \frac{1}{n-1} [x_{j1}^2 + \dots + x_{jn}^2] = \\ &= \frac{1}{n-1} (x_{j1} \ \dots \ x_{jn}) \cdot \begin{pmatrix} x_{j1} \\ \vdots \\ x_{jn} \end{pmatrix} = \\ &= \frac{1}{n-1} \mathbf{x}_j \cdot \mathbf{x}_j\end{aligned}$$

where  $\mathbf{x}_j$  represents the **rows** of  $\mathbf{X}$ !

The total variance over all rows  $j$  reflects the total squared norm

$$\sum_j \sigma_j^2 = \frac{1}{n-1} \sum_j \mathbf{x}_j \cdot \mathbf{x}_j = \frac{1}{n-1} \sum_{j=1}^m \sum_{i=1}^n x_{ji}^2$$

# Covariance matrix

Covariance of dimensions or rows  $j$  and  $k$

$$\sigma_{jk} = \frac{1}{n-1} [x_{j1}x_{k1} + \dots + x_{jn}x_{kn}] = \frac{1}{n} (x_{j1} \dots x_{jn}) \cdot \begin{pmatrix} x_{k1} \\ \vdots \\ x_{kn} \end{pmatrix}$$

$$\mathbf{C}_X = \frac{1}{n-1} \begin{pmatrix} \sigma_1^2 & \sigma_{12} & \dots & \sigma_{1m} \\ \sigma_{21} & \sigma_2^2 & \dots & \sigma_{2m} \\ \vdots & \vdots & \ddots & \vdots \\ \sigma_{m1} & \sigma_{m2} & \dots & \sigma_m^2 \end{pmatrix} =$$

$$= \frac{1}{n-1} \begin{pmatrix} \mathbf{x}_1 \cdot \mathbf{x}_1 & \mathbf{x}_1 \cdot \mathbf{x}_2 & \dots & \mathbf{x}_1 \cdot \mathbf{x}_m \\ \mathbf{x}_2 \cdot \mathbf{x}_1 & \mathbf{x}_2 \cdot \mathbf{x}_2 & \dots & \mathbf{x}_2 \cdot \mathbf{x}_m \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{x}_m \cdot \mathbf{x}_1 & \mathbf{x}_m \cdot \mathbf{x}_2 & \dots & \mathbf{x}_m \cdot \mathbf{x}_m \end{pmatrix} =$$

$$\begin{aligned}
&= \frac{1}{n-1} \begin{pmatrix} \mathbf{x}_1 \cdot \mathbf{x}_1 & \mathbf{x}_1 \cdot \mathbf{x}_2 & \dots & \mathbf{x}_1 \cdot \mathbf{x}_m \\ \mathbf{x}_2 \cdot \mathbf{x}_1 & \mathbf{x}_2 \cdot \mathbf{x}_2 & \dots & \mathbf{x}_2 \cdot \mathbf{x}_m \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{x}_m \cdot \mathbf{x}_1 & \mathbf{x}_m \cdot \mathbf{x}_2 & \dots & \mathbf{x}_m \cdot \mathbf{x}_m \end{pmatrix} = \\
&= \frac{1}{n-1} \begin{pmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \\ \vdots \\ \mathbf{x}_m \end{pmatrix} \begin{pmatrix} \mathbf{x}_1^T & \mathbf{x}_2^T & \dots & \mathbf{x}_m^T \end{pmatrix} = \frac{1}{n-1} \mathbf{X} \mathbf{X}^T \in \mathbb{R}^{m \times m}
\end{aligned}$$

Recall that  $\mathbf{x}_j \in \mathbb{R}^{1 \times n}$ ,  $\mathbf{X} \in \mathbb{R}^{m \times n}$ ,  $\mathbf{X}^T \in \mathbb{R}^{n \times m}$ , and  $\mathbf{C}_\mathbf{X} \in \mathbb{R}^{m \times m}$ .

# Interim summary

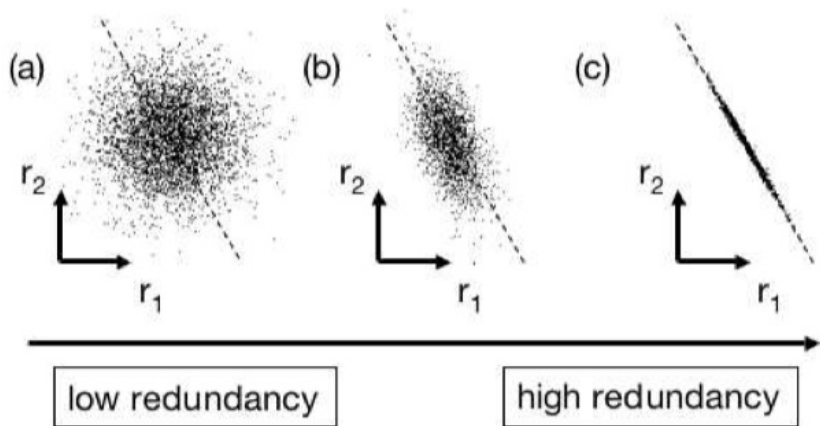
Properties of  $\mathbf{C}_X = \frac{1}{n-1} \mathbf{X}\mathbf{X}^T$

- ▶  $\mathbf{C}_X$  is a square symmetric  $m \times m$  matrix.
- ▶ Its diagonal terms are the *variances* of individual variables (rows of observations  $\mathbf{X}$ ).
- ▶ Its off-diagonal terms are the *covariances* between different variables (two rows of observations  $\mathbf{X}$ ).

$\mathbf{C}_X$  captures the correlations between all pairs of observations.

- ▶ In the diagonal terms, large (small) values indicate signal (noise).
- ▶ In the off-diagonal terms, large (small) values indicate high (low) redundancy.

# Redundancy



# Non-redundant and redundant variance

$$\mathbf{C}_X = \frac{1}{n-1} \underbrace{\begin{pmatrix} \sigma_1^2 & 0 & \dots & 0 \\ 0 & \sigma_2^2 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \sigma_m^2 \end{pmatrix}}_{\text{non-redundant:signal}} +$$
$$+ \frac{1}{n-1} \underbrace{\begin{pmatrix} 0 & \sigma_{12} & \dots & \sigma_{1m} \\ \sigma_{21} & 0 & \dots & \sigma_{2m} \\ \vdots & \vdots & \ddots & \vdots \\ \sigma_{m1} & \sigma_{m2} & \dots & 0 \end{pmatrix}}_{\text{redundant:noise}}$$



# Signal-to-noise ratio

The signal-to-noise ratio is defined as

$$SNR = \frac{\sigma_{signal}^2}{\sigma_{noise}^2}$$

It also specifies the fractional distribution of variance between signal and noise

$$f_{signal} = \frac{SNR}{1 + SNR} = \frac{\sigma_{signal}^2}{\sigma_{signal}^2 + \sigma_{noise}^2}$$

$$f_{noise} = \frac{1}{1 + SNR} = \frac{\sigma_{noise}^2}{\sigma_{signal}^2 + \sigma_{noise}^2}$$

## Choosing transformation $P$

We can now formulate the criteria for choosing a transformation that helps reduce redundancy

$$Y = PX, \quad C_X \rightarrow C_Y$$

- ▶ Seek to maximize variance within as few dimensions of  $Y$  as possible.
- ▶ In identifying the most variable dimensions, we increase signal-to-noise ratio.
- ▶ In identifying redundant dimensions, we discover ways to reduce the dimensionality of our observations.
- ▶ Formally, we seek to maximize a few diagonal terms and minimize the other diagonal and all the off-diagonal terms of  $C_Y$ !
- ▶ A fully diagonal  $C_Y$  is always attainable! How variance distributes over the diagonal depends on our data.

## 4 Diagonalizing variance

There are many ways of choosing  $\mathbf{P}$  such as to diagonalize  $\mathbf{C}_Y$

$$\mathbf{Y} = \mathbf{P}\mathbf{X}, \quad \mathbf{C}_X \rightarrow \mathbf{C}_Y$$

For example:

- ▶ Select the normalized direction in  $\mathbb{R}^{m \times m}$  along which the variance in  $\mathbf{X}$  is maximal. Save this vector as  $\mathbf{p}_1$ .
- ▶ Find a second, orthonormal direction along which the variance  $\mathbf{X}$  is maximal. Save this vector as  $\mathbf{p}_2$ .
- ▶ Repeat until an orthonormal basis with  $m$  vectors has been selected.

Proceeding in this way does not take full benefit of linear algebra!

# Orthogonal diagonalization

Every square symmetric  $\mathbf{S}$  matrix is *orthonormally diagonalizable*

$$\mathbf{S} = \mathbf{E} \mathbf{D} \mathbf{E}^T$$

with square matrix  $\mathbf{E}$  collecting **column** eigenvectors of  $\mathbf{S}$ :

$$\mathbf{E} = (\mathbf{e}_1 \ \mathbf{e}_2 \ \dots \ \mathbf{e}_m), \quad \mathbf{E} \mathbf{E}^T = \mathbf{E}^T \mathbf{E} = \mathbf{I}$$

and diagonal matrix  $\mathbf{D}$  collecting eigenvalues of  $\mathbf{S}$ :

$$\mathbf{D} = \begin{pmatrix} \lambda_1 & \dots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \dots & \lambda_m \end{pmatrix} \quad \text{diag}(\mathbf{D}) = (\lambda_1 \ \lambda_2 \ \dots \ \lambda_m)$$

The relevance of this trick is that it applies to covariance matrices, which are both square and symmetric!

# A matrix is what a matrix does!

$$\mathbf{S} = \underbrace{\mathbf{E}}_{\text{rotate back in } \mathbb{R}^{m \times m}} \underbrace{\mathbf{D}}_{\text{stretch}} \underbrace{\mathbf{E}^T}_{\text{rotate in } \mathbb{R}^{m \times m}} \in \mathbb{R}^{m \times m}$$

After rotation,  $j^{\text{th}}$  component is stretched by  $\lambda_j = D_{jj}$

## Application to covariance

Given orthonormal diagonalization  $\mathbf{C}_X = \mathbf{E} \mathbf{D} \mathbf{E}^T$ , the covariance of  $\mathbf{Y} = \mathbf{P} \mathbf{X}$  may become diagonal:

$$\begin{aligned} \mathbf{C}_Y &= \frac{1}{n-1} \mathbf{Y} \mathbf{Y}^T = \frac{1}{n-1} (\mathbf{P} \mathbf{X})(\mathbf{P} \mathbf{X})^T = \\ &= \frac{1}{n-1} (\mathbf{P} \mathbf{X})(\mathbf{X}^T \mathbf{P}^T) = \frac{1}{n-1} \mathbf{P} (\mathbf{X} \mathbf{X}^T) \mathbf{P}^T = \\ &= \mathbf{P} \mathbf{C}_X \mathbf{P}^T = \mathbf{P} \mathbf{E} \mathbf{D} \mathbf{E}^T \mathbf{P}^T = \mathbf{I} \mathbf{D} \mathbf{I} = \mathbf{D} \end{aligned}$$

if only the **rows** of  $\mathbf{P}$  are chosen as the **columns** of  $\mathbf{E}$ :

$$\mathbf{P} = \mathbf{E}^T \quad \Leftrightarrow \quad \mathbf{P} \mathbf{E} = \mathbf{I} \quad \Leftrightarrow \quad \mathbf{E}^T \mathbf{P}^T = \mathbf{I}$$

## Sort principal components

By convention, we sort rows of  $\mathbf{P}$  (columns of  $\mathbf{E}$ ) from the largest to the smallest eigenvalue (as long as they are non-zero):

The *first principal component*  $\mathbf{p}_1$  captures the largest variance  $\sigma_1^2 = \lambda_1 = D_{1,1}$ .

The *second principal component*  $\mathbf{p}_2$ , captures the next largest variance  $\sigma_2^2 = \lambda_2 = D_{2,2}$ .

⋮

The  *$m^{\text{th}}$  principal component*  $\mathbf{p}_m$ , captures the smallest variance  $\sigma_m^2 = \lambda_m = D_{m,m}$ .

# Different roads to Rome

*Orthonormal diagonalization: requires covariance*

- ▶ Form covariance of observations  $\mathbf{C}_X = \frac{1}{n-1} \mathbf{X} \mathbf{X}^T \in \mathbb{R}^{m \times m}$
- ▶ Orthonormal diagonalization  $\mathbf{C}_X = \mathbf{E} \mathbf{D} \mathbf{E}^T$  into eigenvectors  $\mathbf{E} \in \mathbb{R}^{m \times m}$  and eigenvalues  $\text{diag}(\mathbf{D})$ .
- ▶ Obtain principal components as  $\mathbf{P} = \mathbf{E}^T$  and variance contributions as  $\text{diag}(\mathbf{D})$ .

*Singular-value decomposition (SVD): proceeds directly from observations*

- ▶ Singular value decomposition of  $\mathbf{X}^T = \mathbf{U} \mathbf{\Sigma} \mathbf{V}^T \in \mathbb{R}^{n \times m}$ .
- ▶ Obtain principal components as  $\mathbf{P} = \mathbf{V}^T$  and variance contributions as  $\text{diag}(\mathbf{\Sigma}^T \mathbf{\Sigma})$ .



# Singular Value Decomposition

What the orthonormal diagonalization is for square and symmetric matrices, the *singular value decomposition* is for arbitrary rectangular matrices. Every rectangular matrix  $\mathbf{A}$  has a singular value decomposition

$$\mathbf{A} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T \in \mathbb{R}^{n \times m}$$

with *left singular vectors*

$$\mathbf{U} \in \mathbb{R}^{n \times n} \quad \textit{orthonormal}$$

and *singular values*

$$\mathbf{\Sigma} \in \mathbb{R}^{n \times m} \quad \textit{diagonal}$$

and *right singular vectors*

$$\mathbf{V} \in \mathbb{R}^{m \times m} \quad \textit{orthonormal}$$

Orthonormal (and square)

$$\mathbf{U} = \begin{pmatrix} u_{11} & \dots & u_{1n} \\ \vdots & \ddots & \vdots \\ u_{n1} & \dots & u_{nn} \end{pmatrix} \in \mathbb{R}^{n \times n}, \quad \mathbf{V} = \begin{pmatrix} v_{11} & \dots & v_{1m} \\ \vdots & \ddots & \vdots \\ v_{m1} & \dots & v_{mm} \end{pmatrix} \in \mathbb{R}^{m \times m}$$

Diagonal (and rectangular)

$$\mathbf{\Sigma} = \begin{pmatrix} \sqrt{\lambda_1} & \dots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \dots & \sqrt{\lambda_m} \\ \vdots & \dots & \vdots \\ 0 & \dots & 0 \end{pmatrix} \in \mathbb{R}^{n \times m}$$

# A matrix is what a matrix does!

$$\mathbf{A} = \underbrace{\mathbf{U}}_{\text{rotate in } \mathbb{R}^{n \times n}} \underbrace{\mathbf{\Sigma}}_{\text{stretch and pad}} \underbrace{\mathbf{V}^T}_{\text{rotate in } \mathbb{R}^{m \times m}} \in \mathbb{R}^{n \times m}$$

After first rotation,  $j^{\text{th}}$  component is stretched by  $\sqrt{\lambda_j} = \Sigma_{jj}$

# Application to PCA

Choosing

$$\mathbf{A} = \frac{1}{\sqrt{n-1}} \mathbf{X}^T \in \mathbb{R}^{n \times m}, \quad \text{with} \quad \mathbf{A} = \mathbf{U} \mathbf{\Sigma} \mathbf{V}^T$$

The desired orthonormal basis  $\mathbf{P}$  is

$$\mathbf{P} = \mathbf{V}^T \in \mathbb{R}^{m \times m}$$

and the diagonal matrix of eigenvalues  $\mathbf{D}$  is

$$\mathbf{D} = \begin{pmatrix} \lambda_1 & \dots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \dots & \lambda_m \end{pmatrix} = \mathbf{\Sigma} \mathbf{\Sigma}^T = \begin{pmatrix} \sigma_1^2 & \dots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \dots & \sigma_m^2 \end{pmatrix}$$

# Proof

The choice of  $\mathbf{A} = \frac{1}{\sqrt{n-1}} \mathbf{X}^T = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T$  ensures that

$$\begin{aligned}\mathbf{C}_X &= \frac{1}{n-1} \mathbf{X} \mathbf{X}^T = \mathbf{A}^T \mathbf{A} = \\ &= (\mathbf{U}\mathbf{\Sigma}\mathbf{V}^T)^T (\mathbf{U}\mathbf{\Sigma}\mathbf{V}^T) = \\ &= (\mathbf{V}\mathbf{\Sigma}^T \mathbf{U}^T) (\mathbf{U}\mathbf{\Sigma}\mathbf{V}^T) = \\ &= \mathbf{V}\mathbf{\Sigma}^T \mathbf{\Sigma} \mathbf{V}^T\end{aligned}$$

$$\Leftrightarrow \mathbf{C}_X \mathbf{V} = \mathbf{V}\mathbf{\Sigma}^T \mathbf{\Sigma} = \mathbf{\Sigma}^T \mathbf{\Sigma} \mathbf{V} = \mathbf{D} \mathbf{V}$$

so that  $\mathbf{V}$  are the eigenvectors and  $\mathbf{\Sigma}^T \mathbf{\Sigma}$  the eigenvalues of  $\mathbf{C}_X$ .

## Summary diagonalizing variance

Seek linear transformation  $\mathbf{Y} = \mathbf{P}\mathbf{X}$  such that covariance  $\mathbf{C}_\mathbf{Y}$  is diagonal matrix  $\mathbf{D}$ !

SVD of rectangular matrix  $\mathbf{A} = \frac{1}{n-1}\mathbf{X}^T$  into  $\mathbf{A} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T$  gives us

$$\mathbf{P} = \mathbf{V}^T, \quad \mathbf{D} = \mathbf{\Sigma}^T\mathbf{\Sigma}$$

Orthonormal diagonalization of square and symmetric matrix  $\mathbf{C}_\mathbf{X}$  into  $\mathbf{C}_\mathbf{X} = \mathbf{E}\mathbf{D}\mathbf{E}^T$  gives us

$$\mathbf{P} = \mathbf{E}^T, \quad \mathbf{D} = \mathbf{D}$$

**What have we achieved?** By re-representing observations in PC coordinates, we have eliminated redundancy and revealed main sources of signal (variance):

$$\mathbf{C}_X = \frac{1}{n-1} \underbrace{\begin{pmatrix} \sigma_1^2 & \sigma_{12} & \dots & \sigma_{1m} \\ \sigma_{21} & \sigma_2^2 & \dots & \sigma_{2m} \\ \vdots & \vdots & \ddots & \vdots \\ \sigma_{m1} & \sigma_{m2} & \dots & \sigma_m^2 \end{pmatrix}}_{\text{non-diagonal}}$$

$$\sigma_{total}^2 = \sigma_1^2 + \sigma_2^2 + \dots + \sigma_m^2$$

$$\mathbf{C}_Y = \frac{1}{n-1} \underbrace{\begin{pmatrix} \rho_1^2 & 0 & \dots & 0 \\ 0 & \rho_2^2 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \rho_m^2 \end{pmatrix}}_{\text{diagonal}}$$

$$\sigma_{total}^2 = \rho_1^2 + \rho_2^2 + \dots + \rho_m^2$$

## 5. De-noising with PCA

Full projection matrix  $\mathbf{V}^T \in \mathbb{R}^{m \times m}$

$$\mathbf{Y} = \mathbf{V}^T \mathbf{X}$$

$$\begin{pmatrix} y_{11} & y_{12} & \cdots & y_{1n} \\ \vdots & \ddots & \ddots & \vdots \\ y_{m1} & y_{m2} & \cdots & y_{mn} \end{pmatrix} = \begin{pmatrix} v_{11} & \cdots & v_{m1} \\ \vdots & \ddots & \vdots \\ v_{1m} & \cdots & v_{mm} \end{pmatrix} \begin{pmatrix} x_{11} & x_{1,2} & \cdots & x_{1n} \\ \vdots & \ddots & \ddots & \vdots \\ x_{m1} & x_{m2} & \cdots & x_{mn} \end{pmatrix}$$

Truncated projection matrix  $\tilde{\mathbf{V}}^T \in \mathbb{R}^{r \times m}$ , with  $r < m$

$$\tilde{\mathbf{Y}} = \tilde{\mathbf{V}}^T \mathbf{X}$$

$$\begin{pmatrix} y_{11} & y_{12} & \cdots & y_{1n} \\ y_{r1} & y_{r2} & \cdots & y_{rn} \end{pmatrix} = \begin{pmatrix} v_{11} & \cdots & v_{m1} \\ v_{1r} & \cdots & v_{mr} \end{pmatrix} \begin{pmatrix} x_{11} & x_{1,2} & \cdots & x_{1n} \\ \vdots & \ddots & \vdots & \\ x_{m1} & x_{m2} & \cdots & x_{mn} \end{pmatrix}$$



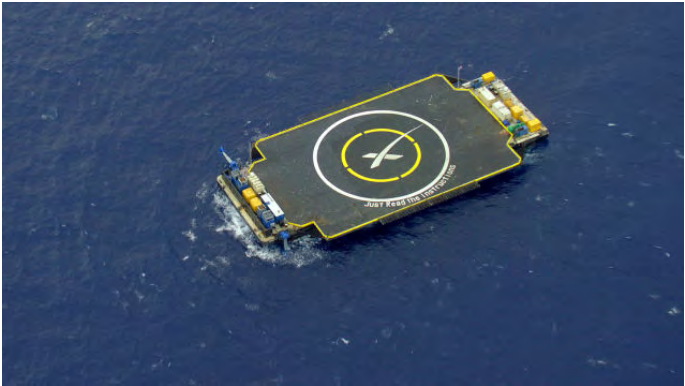
The truncated projection  $\tilde{\mathbf{Y}}$  combines the first  $r$  *principal components* of the original observations  $\mathbf{X}$ . Being smaller, it constitutes a *compressed* version of  $\mathbf{X}$ .

It may be *uncompressed* into a matrix  $\tilde{\mathbf{X}} \in \mathbb{R}^{n \times m}$  by performing the backprojection

$$\tilde{\mathbf{X}} = (\tilde{\mathbf{V}}^T)^{-1} \tilde{\mathbf{Y}}$$

$$\begin{pmatrix} x_{11} & x_{1,2} & \dots & x_{1n} \\ \vdots & \ddots & \ddots & \vdots \\ x_{m1} & x_{m2} & \dots & x_{mn} \end{pmatrix} = \begin{pmatrix} v_{11} & v_{r1} \\ \vdots & \vdots \\ v_{m1} & v_{mr} \end{pmatrix} \begin{pmatrix} y_{11} & y_{12} & \dots & y_{1n} \\ y_{r1} & y_{r2} & \dots & y_{rn} \end{pmatrix}$$

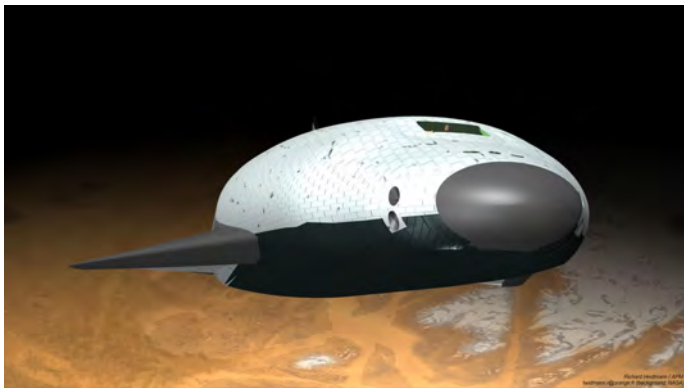
If higher principal components reflect noise, compressing and uncompressing with principal components also serves to *de-noise* the original data.



BW original



BW denoised



BW original

BW denoised

## Summary de-noising

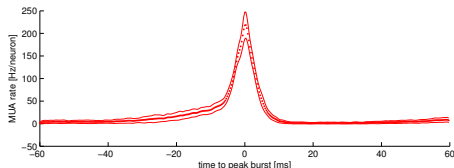
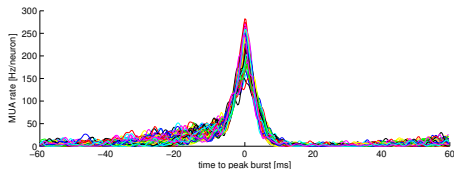
- ▶ Transformed observations  $\mathbf{Y} = \mathbf{P}\mathbf{X}$  span the entire principal component space.
- ▶ It makes sense to consider only the first few dimensions, where most of the variance is concentrated.
- ▶ To this end, transformed observations  $\mathbf{Y}$  may be truncated to  $\tilde{\mathbf{Y}}$  by zeroing the unwanted dimensions.
- ▶ The results may be visualized by back-transforming the truncated observations into the original space

$$\tilde{\mathbf{X}} = \mathbf{P}^T \tilde{\mathbf{Y}}$$

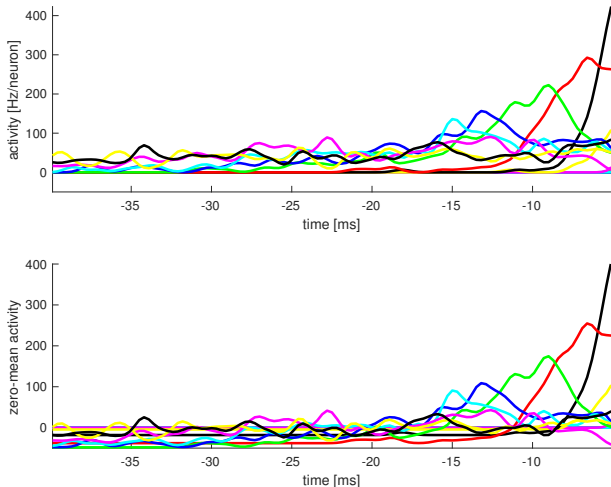
- ▶ The result is a 'de-noised' version of the original observations.

## 6. Multi-unit activity

Consider a network with intermittent population bursts. We have recorded the activity of 400 neurons over approximately 100 bursts, beginning 60 *ms* before and ending 60 *ms* after each burst. We have sorted neurons by overall activity and have saved the average MUA of 20 groups (of 20 neurons each) in 3 *ms*-wide time-bins, spaced 0.25 *ms* apart.

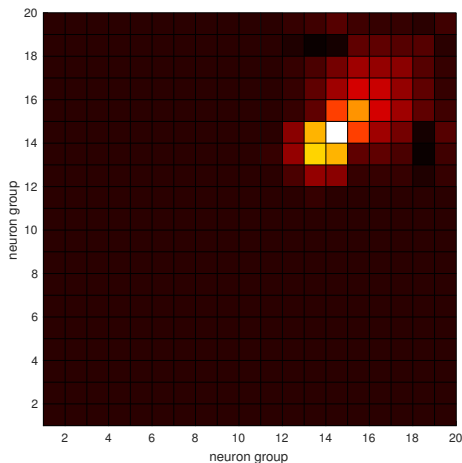


We are interested in the initial stages of each burst and decide to examine the period from 40 *ms* to 5 *ms* before the peak of each burst. Thus we have recordings from approximately 120 time points (spaced 0.25 *ms* apart over 35 *ms*). Notice that neuron groups differ in their level and time of their activity.



## Covariance

Forming activity matrix  $\mathbf{X}_0$  (rows = groups of neurons, columns = time points) and covariance  $\mathbf{C}_X = \mathbf{X}_0 \mathbf{X}_0^T / 120$ , we note that most variance occurs in the more active groups of neurons:

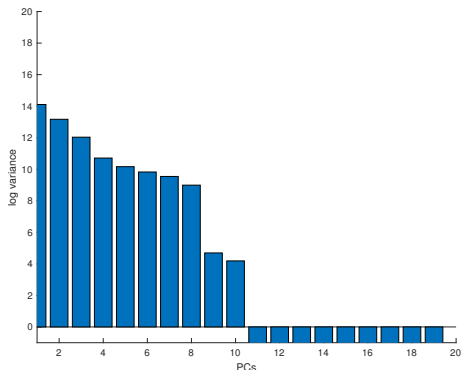


# Variance captured by principal components

Using singular value decomposition

$$\mathbf{X}_0 = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T$$

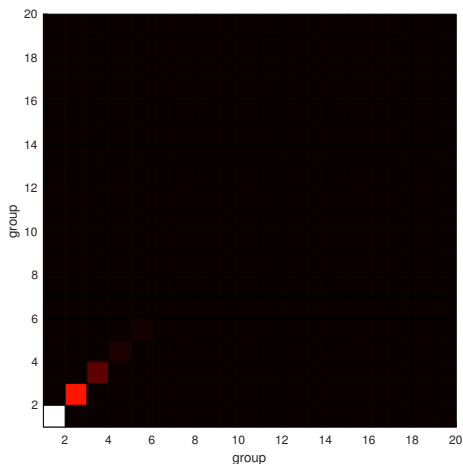
we obtain the principal components  $\mathbf{P} = \mathbf{V}^T$  and the variance  $\text{diag}(\mathbf{\Sigma}\mathbf{\Sigma}^T)$  captured by each:





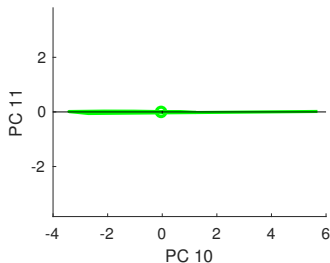
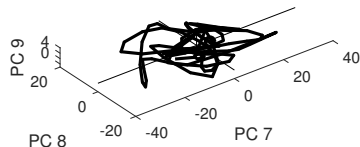
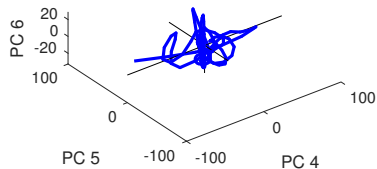
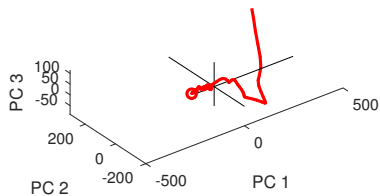
## Covariance of transformed observations

Transforming observations  $\mathbf{Y} = \mathbf{P}\mathbf{X}_0$ , we find that covariance  $\mathbf{C}_Y = \mathbf{Y}\mathbf{Y}^T/120$  is diagonalized, as expected:



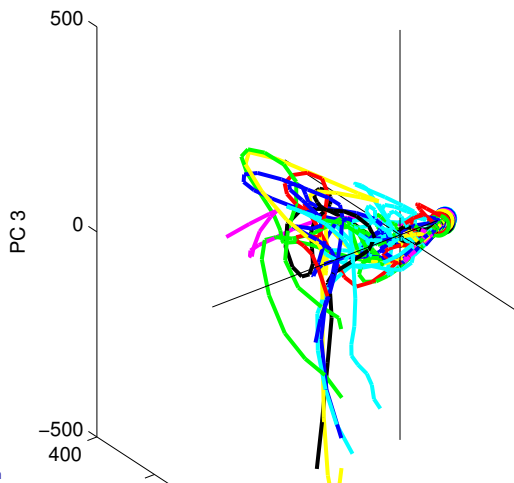
# Visualizing observations in PC space

Having transformed observations, we can visualize the *average burst initiation* in the principal component space:



# Individual burst initiations

In space of first three principal components:

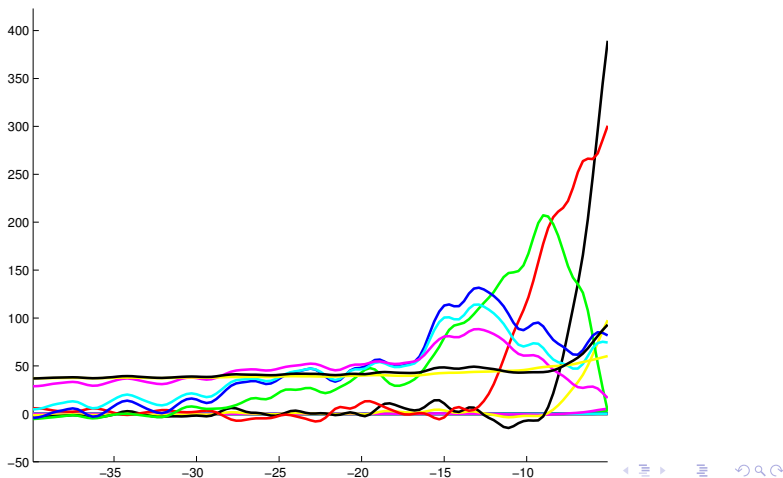


## De-noised observations

Backprojecting the first three rows of  $\mathbf{Y}$  into the original space

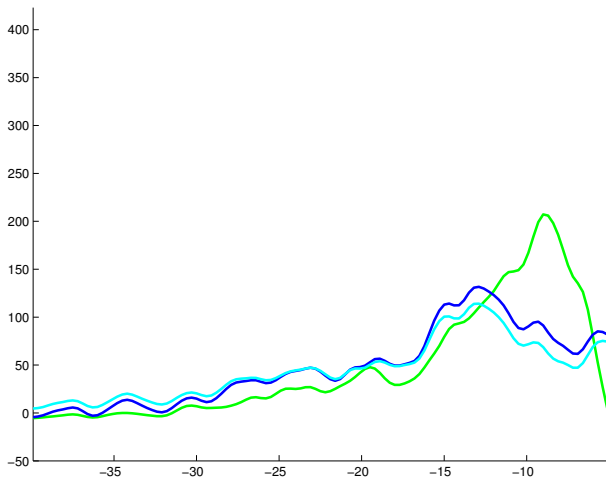
$$\mathbf{X}'_0 = \mathbf{P}^T \mathbf{Y}'$$

we obtain the typical time-course of burst-initiation:



## Potentially important groups of neurons

This draws our attention to three groups of neurons (15, 16, 17) whose activity seems to initiate the burst. To test this hypothesis, we can target these groups with experimental manipulations.



## Summary multi-unit activity

- ▶ We have used PCA to obtain survey a large data set (time-dependent activity of 20 groups of neurons).
- ▶ We find that variance is concentrated in a few groups of neurons.
- ▶ The de-noised time-course reveals the typical development of group activity.
- ▶ This helps identify the groups that initiate bursting activity.

# Summary of PCA assumptions

1. *Linearity*

We frame the problem as a change of basis  $\mathbb{R}^{m \times m}$  space.

2. *Variance is a sufficient statistic*

We assume the variance fully describes the distribution of our (zero-mean) variables.

3. *Large variances are important*

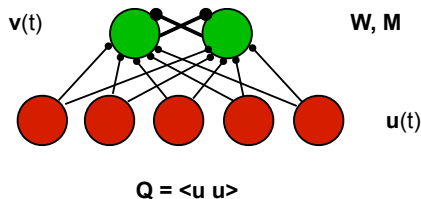
We assume variables with high variance are 'signal', whereas those with low variance are 'noise'.

4. *Principal components are orthogonal*

This lets us perform PCA with linear algebra techniques.

## Link to representational learning

Recall “unsupervised learning” with competing output units.



Learning shaped by

- ▶ *correlations* between input neurons, specifically the eigenvectors of the correlation matrix: *principal components* of input.
- ▶ *competition* between output neurons, so that different neurons can represent different *linear combinations* of principal components.

Output representation in principal component space!



# Next: Sparse Coding

## Appendix: Singular value decomposition

A singular-value decomposition  $\mathbf{A} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T$  decomposes the transformation (change of basis) performed by  $\mathbf{A}$  into three distinct and intuitively intelligible steps. Consider a linear transformation  $\mathbf{x} \rightarrow \mathbf{y}$

$$\mathbf{y} = \mathbf{A}\mathbf{x} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T \mathbf{x}, \quad \mathbf{y} \in \mathbb{R}^n, \quad \mathbf{x} \in \mathbb{R}^m, \quad \mathbf{A} \in \mathbb{R}^{n \times m}$$

Expressing  $\mathbf{y}$  in the orthonormal basis  $\mathbf{U}$  and  $\mathbf{x}$  in the orthonormal basis  $\mathbf{V}$ , we have

$$\hat{\mathbf{y}} = \mathbf{U}\mathbf{y} \quad \Leftrightarrow \quad \mathbf{y} = \mathbf{U}^T \hat{\mathbf{y}}, \quad \hat{\mathbf{x}} = \mathbf{V}\mathbf{x} \quad \Leftrightarrow \quad \mathbf{x} = \mathbf{V}^T \hat{\mathbf{x}}$$

By left-multiplying the SVD with  $\mathbf{U}^T$ , we can reformulate the transformation in terms of  $\hat{\mathbf{x}} \rightarrow \hat{\mathbf{y}}$

$$\mathbf{U}^T \cdot \left| \quad \mathbf{y} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T \mathbf{x} \quad \Leftrightarrow \quad \hat{\mathbf{y}} = \mathbf{U}^T \mathbf{y} = \mathbf{\Sigma} \hat{\mathbf{x}} \right.$$

Thus, in terms of the orthonormal bases  $\mathbf{U}$  and  $\mathbf{V}$ , the transformation *rescales* components 1 to  $m$  with the singular values in  $\mathbf{\Sigma}$  and zeroes all others:

$$\begin{pmatrix} \hat{y}_1 \\ \vdots \\ \hat{x}_m \\ \vdots \\ 0 \end{pmatrix} = \begin{pmatrix} \sqrt{\lambda_1} & \dots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \dots & \sqrt{\lambda_m} \\ \vdots & \dots & \vdots \\ 0 & \dots & 0 \end{pmatrix} \begin{pmatrix} \hat{x}_1 \\ \hat{x}_2 \\ \vdots \\ \hat{x}_m \end{pmatrix}$$

Note that  $\mathbf{\Sigma}^T \mathbf{\Sigma}$  is square and diagonal:

$$\mathbf{\Sigma}^T \mathbf{\Sigma} = \mathbf{\Sigma} \mathbf{\Sigma}^T = \begin{pmatrix} \lambda_1 & \dots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \dots & \lambda_m \end{pmatrix} \in \mathbb{R}^{m \times m}$$

# Derivation of SVD

Consider a covariance matrix and its orthonormal eigenvectors and real positive eigenvalues

$$\mathbf{X}^T \mathbf{X} \mathbf{v}_j = \lambda_j \mathbf{v}_j, \quad \mathbf{X} \in \mathbb{R}^{n \times m}, \mathbf{v}_j \in \mathbb{R}^{m \times 1},$$

Use  $\mathbf{X}$  to transform  $m \times 1$  eigenvectors  $\mathbf{v}_j$  into  $n \times 1$  vectors  $\mathbf{u}_j$

$$\mathbf{u}_j = \frac{1}{\sqrt{\lambda_j}} \mathbf{X} \mathbf{v}_j, \quad \mathbf{u}_j \in \mathbb{R}^{n \times 1}$$

Astonishingly, the transformed vectors also form an orthonormal basis! Their norm is unity! They are pairwise orthogonal!

$$\|\mathbf{u}_j\| = 1, \quad \mathbf{u}_j \mathbf{u}_k^T = \delta_{jk}$$

# Proof

$$\begin{aligned}\|\mathbf{u}_j\|^2 &= \mathbf{u}_j \cdot \mathbf{u}_j = \mathbf{u}_j^T \mathbf{u}_j = \\ &= \frac{1}{\lambda_j} (\mathbf{X} \mathbf{v}_j)^T (\mathbf{X} \mathbf{v}_j) = \frac{1}{\lambda_j} (\mathbf{v}_j^T \mathbf{X}^T) (\mathbf{X} \mathbf{v}_j) = \\ &= \frac{1}{\lambda_j} \mathbf{v}_j^T \lambda_j \mathbf{v}_j = \mathbf{v}_j \cdot \mathbf{v}_j = \|\mathbf{v}_j\|^2 = 1\end{aligned}$$

$$\begin{aligned}\mathbf{u}_j \cdot \mathbf{u}_k &= \mathbf{u}_j^T \mathbf{u}_k = \\ &= \frac{1}{\sqrt{\lambda_j \lambda_k}} (\mathbf{X} \mathbf{v}_j)^T (\mathbf{X} \mathbf{v}_k) = \frac{1}{\sqrt{\lambda_j \lambda_k}} (\mathbf{v}_j^T \mathbf{X}^T) (\mathbf{X} \mathbf{v}_k) = \\ &= \frac{1}{\sqrt{\lambda_j \lambda_k}} \mathbf{v}_j^T \lambda_k \mathbf{v}_k = \sqrt{\frac{\lambda_k}{\lambda_j}} \mathbf{v}_j \cdot \mathbf{v}_k = \delta_{jk}\end{aligned}$$

## Collect into matrix

Collect  $\mathbf{u}$ 's into  $n \times m$  matrix

$$\mathbf{U}' = (\mathbf{u}_1 \dots \mathbf{u}_m) \in \mathbb{R}^{n \times m}$$

'Pad' with further orthonormal vectors to complete square matrix

$$\mathbf{U} = (\mathbf{u}_1 \quad \dots \quad \mathbf{u}_m \quad \mathbf{u}'_{m+1} \quad \dots \quad \mathbf{u}'_n) \in \mathbb{R}^{n \times n}$$

Now have orthonormal basis

$$\mathbf{U}\mathbf{U}^T = \mathbf{U}^T\mathbf{U} = \mathbf{I} \quad \in \mathbb{R}^{n \times n}$$

ctd

Collect the individual vector transformation

$$\mathbf{u}_j \sqrt{\lambda_j} = \mathbf{X} \mathbf{v}_j$$

into matrix transformation

$$\mathbf{U} \mathbf{\Sigma} = \mathbf{X} \mathbf{V} \quad \in R^{n \times m}$$

$$\mathbf{U} \in R^{n \times n}, \quad \mathbf{\Sigma} \in R^{n \times m}, \quad \mathbf{X} \in R^{n \times m}, \quad \mathbf{V} \in R^{m \times m},$$

$$\mathbf{\Sigma} = \begin{pmatrix} \sqrt{\lambda_1} & 0 & \dots & 0 \\ 0 & \sqrt{\lambda_2} & \ddots & \vdots \\ \vdots & \ddots & \ddots & \vdots \\ 0 & 0 & \dots & \sqrt{\lambda_m} \\ \vdots & \vdots & \vdots & \vdots \\ 0 & 0 & 0 & 0 \end{pmatrix} \in R^{n \times m}$$

Finally, right-multiply with  $\mathbf{V}^T$  to obtain

$$\mathbf{U}\boldsymbol{\Sigma} = \mathbf{X}\mathbf{V} \quad | \cdot \mathbf{V}^T$$

$$\mathbf{U}\boldsymbol{\Sigma}\mathbf{V}^T = \mathbf{X}$$

To recapitulate, for a rectangular matrix  $\mathbf{X} \in R^{n \times m}$ , we considered the eigenvectors and eigenvalues of the covariance matrix

$$\mathbf{X}^T \mathbf{X} \mathbf{v}_j = \lambda_j \mathbf{v}_j$$

and obtained the singular value decomposition of our matrix  $\mathbf{X}$

$$\mathbf{U}\boldsymbol{\Sigma}\mathbf{V}^T = \mathbf{X}, \quad \Sigma_{jj} = \sqrt{\lambda_j}, \quad \mathbf{V} = ( \mathbf{v}_1 \quad \dots \quad \mathbf{v}_m )$$