

# Principal Component Analysis (PCA)

# Dimensionality reduction

One of the reasons the SVD is used is for dimensionality reduction. However, SVD has many many other uses

Now we look at another dimensionality reduction technique, PCA

PCA is often used as a blackbox technique, here we take a look at the math behind it

# What is PCA?

Linear algebraic technique

Helps reduce a complex dataset to a lower dimensional one

Non-parametric method: does not assume anything about data distribution (distribution from the statistical point of view)

## Brief “review” of some probability concepts

Proper definition of *probability* requires to use *measure theory*..  
will not get into details here

A **random variable**  $X$  is a *measurable* function  $X : \Omega \rightarrow E$ , where  $\Omega$  is a set of outcomes (*sample space*) and  $E$  is a measurable space

$$\mathbb{P}(X \in S \subseteq E) = \mathbb{P}(\omega \in \Omega | X(\omega) \in S)$$

**Distribution function** of a r.v.,  $F(x) = \mathbb{P}(X \leq x)$ , describes the distribution of a r.v.

R.v. can be discrete or continuous or .. other things.

### Definition 1 (Variance)

Let  $X$  be a random variable. The **variance** of  $X$  is given by

$$\text{Var } X = E \left[ (X - E(X))^2 \right]$$

where  $E$  is the expected value

### Definition 2 (Covariance)

Let  $X, Y$  be jointly distributed random variables. The **covariance** of  $X$  and  $Y$  is given by

$$\text{cov}(X, Y) = E [(X - E(X)) (Y - E(Y))]$$

Note that  $\text{cov}(X, X) = E \left[ (X - E(X))^2 \right] = \text{Var } X$

## In practice: “true law” versus “observation”

In statistics: we reason on the *true law* of distributions, but we usually have only access to a sample

We then use **estimators** to .. estimate the value of a parameter, e.g., the mean, variance and covariance

### Definition 3 (Unbiased estimators of the mean and variance)

Let  $x_1, \dots, x_n$  be data points (the *sample*) and

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$$

be the **mean** of the data. An unbiased estimator of the variance of the sample is

$$\sigma^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2$$

#### Definition 4 (Unbiased estimator of the covariance)

Let  $(x_1, y_1), \dots, (x_n, y_n)$  be data points,

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i \text{ and } \bar{y} = \frac{1}{n} \sum_{i=1}^n y_i$$

be the means of the data. An estimator of the covariance of the sample is

$$\text{cov}(x, y) = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})$$



# What does covariance do?

Variance explains how data disperses around the mean, in a 1-D case

Covariance measures the relationship between two dimensions.  
E.g., height and weight

More than the exact value, the sign is important:

- ▶  $\text{cov}(X, Y) > 0$ : both dimensions change in the same “direction”; e.g., larger height usually means higher weight
- ▶  $\text{cov}(X, Y) < 0$ : both dimensions change in reverse directions; e.g., time spent on social media and performance in this class
- ▶  $\text{cov}(X, Y) = 0$ : the dimensions are independent from one another; e.g., sex/gender and “intelligence”

# The covariance matrix

Typically, we consider more than 2 variables..

## Definition 5

Suppose  $p$  random variables  $X_1, \dots, X_p$ . Then the covariance matrix is the symmetric matrix

$$\begin{pmatrix} \text{cov}(X_1, X_1) & \text{cov}(X_1, X_2) & \cdots & \text{cov}(X_1, X_p) \\ \text{cov}(X_2, X_1) & \text{cov}(X_2, X_2) & \cdots & \text{cov}(X_2, X_p) \\ \vdots & \vdots & & \vdots \\ \text{cov}(X_p, X_1) & \text{cov}(X_p, X_2) & \cdots & \text{cov}(X_p, X_p) \end{pmatrix}$$

i.e., using the properties of covariance,

$$\begin{pmatrix} \text{Var } X_1 & \text{cov}(X_1, X_2) & \cdots & \text{cov}(X_1, X_p) \\ \text{cov}(X_1, X_2) & \text{Var } X_2 & \cdots & \text{cov}(X_2, X_p) \\ \vdots & \vdots & & \vdots \\ \text{cov}(X_1, X_p) & \text{cov}(X_2, X_p) & \cdots & \text{Var } X_p \end{pmatrix}$$

## Example of a PCA problem

We collect a bunch of information about a bunch of people.. for instance this data from Loughborough University

*This dataset contains the height, weight and 4 fingerprint measurements (length, width, area and circumference), collected from 200 participants.*

What best describes a participant?

# The variables

Each participant is associated to 11 variables

- ▶ "Participant Number"
- ▶ "Gender"
- ▶ "Age"
- ▶ "Dominant Hand"
- ▶ "Height (cm) (average of 3 measurements)"
- ▶ "Weight (kg) (average of 3 measurements)"
- ▶ "Fingertip Temperature (°C)"
- ▶ "Fingerprint Height (mm)"
- ▶ "Fingerprint Width (mm)"
- ▶ "Fingerprint Area (mm<sup>2</sup>)"
- ▶ "Fingerprint Circumference (mm)"

# Nature of variables

Variables have different natures

- ▶ "Participant Number":  $\in \mathbb{N}$  (not interesting)
- ▶ "Gender": categorical
- ▶ "Age":  $\in \mathbb{N}$
- ▶ "Dominant Hand": categorical
- ▶ "Height (cm) (average of 3 measurements)":  $\in \mathbb{R}$
- ▶ "Weight (kg) (average of 3 measurements)":  $\in \mathbb{R}$
- ▶ "Fingertip Temperature ( $^{\circ}\text{C}$ )":  $\in \mathbb{R}$
- ▶ "Fingerprint Height (mm)":  $\in \mathbb{R}$
- ▶ "Fingerprint Width (mm)":  $\in \mathbb{R}$
- ▶ "Fingerprint Area ( $\text{mm}^2$ )":  $\in \mathbb{R}$
- ▶ "Fingerprint Circumference (mm)":  $\in \mathbb{R}$

## Setting things up

Each participant is a row in the matrix (an *observation*)

Each variable is a column

So we have an  $200 \times 10$  matrix (we discard the “Participant number” column)

We want to find what carries the most information

For this, we are going to project the information in a new basis in which the first “dimension” will carry most variance, the second dimension will carry a little less, etc.

In order to do so, we need to learn how to change bases

# Change of basis

## Definition 6 (Change of basis matrix)

$\mathcal{B} = \{\mathbf{u}_1, \dots, \mathbf{u}_n\}$  and  $\mathcal{C} = \{\mathbf{v}_1, \dots, \mathbf{v}_n\}$  bases of vector space  $V$

The **change of basis matrix**  $P_{\mathcal{C} \leftarrow \mathcal{B}} \in \mathcal{M}_n$ ,

$$P_{\mathcal{C} \leftarrow \mathcal{B}} = [[\mathbf{u}_1]_{\mathcal{C}} \cdots [\mathbf{u}_n]_{\mathcal{C}}]$$

has columns the coordinate vectors  $[\mathbf{u}_1]_{\mathcal{C}}, \dots, [\mathbf{u}_n]_{\mathcal{C}}$  of the vectors in  $\mathcal{B}$  with respect to  $\mathcal{C}$

## Theorem 7

$\mathcal{B} = \{\mathbf{u}_1, \dots, \mathbf{u}_n\}$  and  $\mathcal{C} = \{\mathbf{v}_1, \dots, \mathbf{v}_n\}$  bases of vector space  $V$  and  $P_{\mathcal{C} \leftarrow \mathcal{B}}$  a change of basis matrix from  $\mathcal{B}$  to  $\mathcal{C}$

1.  $\forall \mathbf{x} \in V, P_{\mathcal{C} \leftarrow \mathcal{B}}[\mathbf{x}]_{\mathcal{B}} = [\mathbf{x}]_{\mathcal{C}}$
2.  $P_{\mathcal{C} \leftarrow \mathcal{B}}$  s.t.  $\forall \mathbf{x} \in V, P_{\mathcal{C} \leftarrow \mathcal{B}}[\mathbf{x}]_{\mathcal{B}} = [\mathbf{x}]_{\mathcal{C}}$  is **unique**
3.  $P_{\mathcal{C} \leftarrow \mathcal{B}}$  invertible and  $P_{\mathcal{C} \leftarrow \mathcal{B}}^{-1} = P_{\mathcal{B} \leftarrow \mathcal{C}}$

# Row-reduction method for changing bases

## Theorem 8

$\mathcal{B} = \{\mathbf{u}_1, \dots, \mathbf{u}_n\}$  and  $\mathcal{C} = \{\mathbf{v}_1, \dots, \mathbf{v}_n\}$  bases of vector space  $V$ .  
Let  $\mathcal{E}$  be any basis for  $V$ ,

$$B = [[\mathbf{u}_1]_{\mathcal{E}}, \dots, [\mathbf{u}_n]_{\mathcal{E}}] \text{ and } C = [[\mathbf{v}_1]_{\mathcal{E}}, \dots, [\mathbf{v}_n]_{\mathcal{E}}]$$

and let  $[C|B]$  be the augmented matrix constructed using  $C$  and  $B$ . Then

$$\text{RREF}([C|B]) = [\mathbb{I} | P_{\mathcal{C} \leftarrow \mathcal{B}}]$$

If working in  $\mathbb{R}^n$ , this is quite useful with  $\mathcal{E}$  the standard basis of  $\mathbb{R}^n$  (it does not matter if  $\mathcal{B} = \mathcal{E}$ )



So the question now becomes

*How do we find what new basis to look at our data in?*

(Changing the basis does not change the data, just the view you have of it)

(Think of what happens when you do a headstand.. your up becomes down, your right and left switch, but the world does not change, just your view of it)

(Changes of bases are *fundamental* operations in Science)

## Setting things up

I will use notation (mostly) as in Joliffe's *Principal Component Analysis* (PDF of older version available for free from UofM Libraries)

$\mathbf{x} = (x_1, \dots, x_p)$  vector of  $p$  random variables

We seek a linear function  $\alpha_1^T \mathbf{x}$  with maximum variance, where  $\alpha_1 = (\alpha_{11}, \dots, \alpha_{1p})$ , i.e.,

$$\alpha_1^T \mathbf{x} = \sum_{j=1}^p \alpha_{1j} x_j$$

Then we seek a linear function  $\alpha_2^T \mathbf{x}$  with maximum variance, uncorrelated to  $\alpha_1^T \mathbf{x}$

And we continue...

At  $k$ th stage, we find a linear function  $\alpha_k^T \mathbf{x}$  with maximum variance, uncorrelated to  $\alpha_1^T \mathbf{x}, \dots, \alpha_{k-1}^T \mathbf{x}$

$\alpha_i^T \mathbf{x}$  is the  $i$ th **principal component** (PC)

## Case of known covariance matrix

Suppose we know  $\Sigma$ , covariance matrix of  $\mathbf{x}$  (i.e., typically: we know  $\mathbf{x}$ )

Then the  $k$ th PC is

$$z_k = \alpha_k^T \mathbf{x}$$

where  $\alpha_k$  is an eigenvector of  $\Sigma$  corresponding to the  $k$ th largest eigenvalue  $\lambda_k$

If, additionally,  $\|\alpha_k\| = \alpha_k^T \alpha_k = 1$ , then  $\lambda_k = \text{Var } z_k$

## Why is that?

Let us start with

$$\alpha_1^T \mathbf{x}$$

We want maximum variance, where  $\alpha_1 = (\alpha_{11}, \dots, \alpha_{1p})$ , i.e.,

$$\alpha_1^T \mathbf{x} = \sum_{j=1}^p \alpha_{1j} x_j$$

with the constraint that  $\|\alpha_1\| = 1$

We have

$$\text{Var } \alpha_1^T \mathbf{x} = \alpha_1^T \Sigma \alpha_1$$

# Objective

We want to maximise  $\text{Var } \alpha_1^T \mathbf{x}$ , i.e.,

$$\alpha_1^T \Sigma \alpha_1$$

under the constraint that  $\|\alpha_1\| = 1$

$\implies$  use **Lagrange multipliers**

# Maximisation using Lagrange multipliers

(A.k.a. super-brief intro to multivariable calculus)

We want the max of  $f(x_1, \dots, x_n)$  under the constraint  
 $g(x_1, \dots, x_n) = k$

1. Solve

$$\begin{aligned}\nabla f(x_1, \dots, x_n) &= \lambda \nabla g(x_1, \dots, x_n) \\ g(x_1, \dots, x_n) &= k\end{aligned}$$

where  $\nabla = (\frac{\partial}{\partial x_1}, \dots, \frac{\partial}{\partial x_n})$  is the **gradient operator**

2. Plug all solutions into  $f(x_1, \dots, x_n)$  and find maximum values  
(provided values exist and  $\nabla g \neq \mathbf{0}$  there)

$\lambda$  is the **Lagrange multiplier**

# The gradient

(Continuing our super-brief intro to multivariable calculus)

$f : \mathbb{R}^n \rightarrow \mathbb{R}$  function of several variables,  $\nabla = \left( \frac{\partial}{\partial x_1}, \dots, \frac{\partial}{\partial x_n} \right)$  the gradient operator

Then

$$\nabla f = \left( \frac{\partial}{\partial x_1} f, \dots, \frac{\partial}{\partial x_n} f \right)$$

So  $\nabla f$  is a *vector-valued* function,  $\nabla f : \mathbb{R}^n \rightarrow \mathbb{R}^n$ ; also written as

$$\nabla f = f_{x_1}(x_1, \dots, x_n) \mathbf{e}_1 + \dots + f_{x_n}(x_1, \dots, x_n) \mathbf{e}_n$$

where  $f_{x_i}$  is the partial derivative of  $f$  with respect to  $x_i$  and  $\{\mathbf{e}_1, \dots, \mathbf{e}_n\}$  is the standard basis of  $\mathbb{R}^n$



# Bear with me..

(You may experience a brief period of discomfort)

$\alpha_1^T \Sigma \alpha_1$  and  $\|\alpha_1\|^2 = \alpha_1^T \alpha_1$  are functions of  $\alpha_1 = (\alpha_{11}, \dots, \alpha_{1p})$

In the notation of the previous slide, we want the max of

$$f(\alpha_{11}, \dots, \alpha_{1p}) := \alpha_1^T \Sigma \alpha_1$$

under the constraint that

$$g(\alpha_{11}, \dots, \alpha_{1p}) := \alpha_1^T \alpha_1 = 1$$

and with gradient operator

$$\nabla = \left( \frac{\partial}{\partial \alpha_{11}}, \dots, \frac{\partial}{\partial \alpha_{1p}} \right)$$

## Effect of $\nabla$ on $g$

$g$  is easiest to see:

$$\begin{aligned}\nabla g(\alpha_{11}, \dots, \alpha_{1p}) &= \left( \frac{\partial}{\partial \alpha_{11}}, \dots, \frac{\partial}{\partial \alpha_{1p}} \right) (\alpha_{11}, \dots, \alpha_{1p}) \begin{pmatrix} \alpha_{11} \\ \vdots \\ \alpha_{1p} \end{pmatrix} \\ &= \left( \frac{\partial}{\partial \alpha_{11}}, \dots, \frac{\partial}{\partial \alpha_{1p}} \right) (\alpha_{11}^2 + \dots + \alpha_{1p}^2) \\ &= (2\alpha_{11}, \dots, 2\alpha_{1p}) \\ &= 2\alpha_1\end{aligned}$$

(And that's a general result:  $\nabla \|\mathbf{x}\|_2^2 = 2\mathbf{x}$  with  $\|\cdot\|_2$  the Euclidean norm)

## Effect of $\nabla$ on $f$

Expand (write  $\Sigma = [s_{ij}]$  and do not exploit symmetry)

$$\begin{aligned}\alpha_1^T \Sigma \alpha_1 &= (\alpha_{11}, \dots, \alpha_{1p}) \begin{pmatrix} s_{11} & s_{12} & \cdots & s_{1p} \\ s_{21} & s_{22} & \cdots & s_{2p} \\ \vdots & \vdots & & \vdots \\ s_{p1} & s_{p2} & & s_{pp} \end{pmatrix} \begin{pmatrix} \alpha_{11} \\ \alpha_{12} \\ \vdots \\ \alpha_{1p} \end{pmatrix} \\ &= (\alpha_{11}, \dots, \alpha_{1p}) \begin{pmatrix} s_{11}\alpha_{11} + s_{12}\alpha_{12} + \cdots + s_{1p}\alpha_{1p} \\ s_{21}\alpha_{11} + s_{22}\alpha_{12} + \cdots + s_{2p}\alpha_{1p} \\ \vdots \\ s_{p1}\alpha_{11} + s_{p2}\alpha_{12} + \cdots + s_{pp}\alpha_{1p} \end{pmatrix} \\ &= (s_{11}\alpha_{11} + s_{12}\alpha_{12} + \cdots + s_{1p}\alpha_{1p})\alpha_{11} \\ &\quad + (s_{21}\alpha_{11} + s_{22}\alpha_{12} + \cdots + s_{2p}\alpha_{1p})\alpha_{12} \\ &\quad \vdots \\ &\quad + (s_{p1}\alpha_{11} + s_{p2}\alpha_{12} + \cdots + s_{pp}\alpha_{1p})\alpha_{1p}\end{aligned}$$

We have

$$\begin{aligned}
 \alpha_1^T \Sigma \alpha_1 &= (s_{11}\alpha_{11} + s_{12}\alpha_{12} + \cdots + s_{1p}\alpha_{1p})\alpha_{11} \\
 &\quad + (s_{21}\alpha_{11} + s_{22}\alpha_{12} + \cdots + s_{2p}\alpha_{1p})\alpha_{12} \\
 &\quad \vdots \\
 &\quad + (s_{p1}\alpha_{11} + s_{p2}\alpha_{12} + \cdots + s_{pp}\alpha_{1p})\alpha_{1p}
 \end{aligned}$$

So

$$\begin{aligned}
 \frac{\partial}{\partial \alpha_{11}} \alpha_1^T \Sigma \alpha_1 &= (s_{11}\alpha_{11} + s_{12}\alpha_{12} + \cdots + s_{1p}\alpha_{1p}) + s_{11}\alpha_{11} \\
 &\quad + s_{21}\alpha_{12} \\
 &\quad \vdots \\
 &\quad + s_{p1}\alpha_{1p} \\
 &= s_{11}\alpha_{11} + s_{12}\alpha_{12} + \cdots + s_{1p}\alpha_{1p} \\
 &\quad + s_{11}\alpha_{11} + s_{21}\alpha_{12} + \cdots + s_{p1}\alpha_{1p} \\
 &= 2(s_{11}\alpha_{11} + s_{12}\alpha_{12} + \cdots + s_{1p}\alpha_{1p})
 \end{aligned}$$

(last equality stems from symmetry of  $\Sigma$ )

In general, for  $i = 1, \dots, p$ ,

$$\begin{aligned}\frac{\partial}{\partial \alpha_{1i}} \alpha_1^T \Sigma \alpha_1 &= s_{i1} \alpha_{11} + s_{i2} \alpha_{12} + \dots + s_{ip} \alpha_{1p} \\ &\quad + s_{i1} \alpha_{11} + s_{i2} \alpha_{12} + \dots + s_{ip} \alpha_{1p} \\ &= 2(s_{i1} \alpha_{11} + s_{i2} \alpha_{12} + \dots + s_{ip} \alpha_{1p})\end{aligned}$$

(because of symmetry of  $\Sigma$ )

As a consequence,

$$\nabla \alpha_1^T \Sigma \alpha_1 = 2 \Sigma \alpha_1$$

So solving

$$\nabla f(x_1, \dots, x_n) = \lambda \nabla g(x_1, \dots, x_n)$$

means solving

$$2\Sigma\alpha_1 = \lambda 2\alpha_1$$

i.e.,

$$\Sigma\alpha_1 = \lambda\alpha_1$$

$\implies (\lambda, \alpha_1)$  eigenpair of  $\Sigma$ , with  $\alpha_1$  having unit length

## Picking the right eigenvalue

$(\lambda, \alpha_1)$  eigenpair of  $\Sigma$ , with  $\alpha_1$  having unit length

But which  $\lambda$  to choose?

Recall that we want  $\text{Var } \alpha_1^T \mathbf{x} = \alpha_1^T \Sigma \alpha_1$  maximal

We have

$$\text{Var } \alpha_1^T \mathbf{x} = \alpha_1^T \Sigma \alpha_1 = \alpha_1^T (\Sigma \alpha_1) = \alpha_1^T (\lambda \alpha_1) = \lambda (\alpha_1^T \alpha_1) = \lambda$$

$\implies$  we pick  $\lambda = \lambda_1$ , the largest eigenvalue (covariance matrix symmetric so eigenvalues real)

## What we have this far..

The first principal component is  $\alpha_1^T \mathbf{x}$  and has variance  $\lambda_1$ , where  $\lambda_1$  the largest eigenvalue of  $\Sigma$  and  $\alpha_1$  an associated eigenvector with  $\|\alpha_1\| = 1$

We want the second principal component to be *uncorrelated* with  $\alpha_1^T \mathbf{x}$  and to have maximum variance  $\text{Var } \alpha_2^T \mathbf{x} = \alpha_2^T \Sigma \alpha_2$ , under the constraint that  $\|\alpha_2\| = 1$

$\alpha_2^T \mathbf{x}$  uncorrelated to  $\alpha_1^T \mathbf{x}$  if  $\text{cov}(\alpha_1^T \mathbf{x}, \alpha_2^T \mathbf{x}) = 0$



We have

$$\begin{aligned}\text{cov}(\alpha_1^T \mathbf{x}, \alpha_2^T \mathbf{x}) &= \alpha_1^T \Sigma \alpha_2 \\ &= \alpha_2^T \Sigma^T \alpha_1 \\ &= \alpha_2^T \Sigma \alpha_1 \quad [\Sigma \text{ symmetric}] \\ &= \alpha_2^T (\lambda_1 \alpha_1) \\ &= \lambda \alpha_2^T \alpha_1\end{aligned}$$

So  $\alpha_2^T \mathbf{x}$  uncorrelated to  $\alpha_1^T \mathbf{x}$  if  $\alpha_1 \perp \alpha_2$

This is beginning to sound a lot like Gram-Schmidt, no?

## In short

Take whatever covariance matrix is available to you (known  $\Sigma$  or sample  $S_X$ ) – assume sample from now on for simplicity

For  $i = 1, \dots, p$ , the  $i$ th principal component is

$$z_i = \mathbf{v}_i^T \mathbf{x}$$

where  $\mathbf{v}_i$  eigenvector of  $S_X$  associated to the  $i$ th largest eigenvalue  $\lambda_i$

If  $\mathbf{v}_i$  is normalised, then  $\lambda_i = \text{Var } z_k$

## Covariance matrix

$\Sigma$  the covariance matrix of the random variable,  $S_X$  the sample covariance matrix

$X \in \mathcal{M}_{mp}$  the data, then the (sample) covariance matrix  $S_X$  takes the form

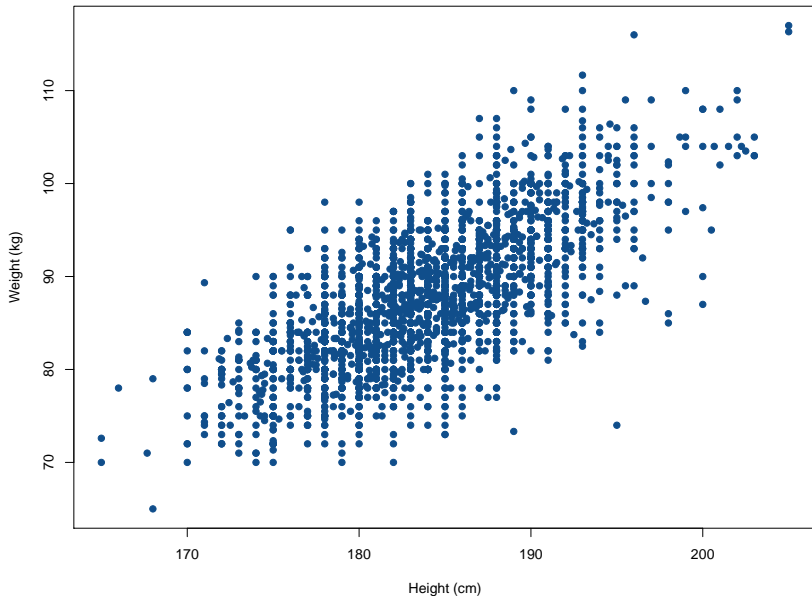
$$S_X = \frac{1}{n-1} X^T X$$

where the data is centred!

Sometimes you will see  $S_X = 1/(n-1)XX^T$ . This is for matrices with observations in columns and variables in rows. Just remember that you want the covariance matrix to have size the number of variables, not observations, this will give you the order in which to take the product

# A smaller 2D example

Hockey players at IIHF world championships 2001-2016



# Centre the data

Subtract the mean (our first – simple – change of basis)

