

---

---

---

---

---



# Week 1

High dimensional data: difficult to interpret, hard to analyse, difficult to visualise, storage can be expensive

↳ often overcomplete: many dimensions are redundant, can be explained by combination of other dimensions

Dimensionality reduction exploits structure and correlation

↳ work with more compact representation of the data, without losing info

↳ can think of as a compression technique (e.g. JPEG)

e.g. digit "8": 28x28 pixels, 784D vector (pixels are structured, not random)

↙ "feature"

↳ many examples of 8's: differ slightly, can use D.R. to find lower dimensional representation of all 8's (easier to work with)

PCA · linear D.R.

## Mean Values

Use statistical properties to describe the data

Mean = average data point, not necessarily part of dataset (e.g. average "8" image)

$$D = \{x_1, \dots, x_N\}$$
$$E[D] = \frac{1}{N} \sum_{n=1}^N x_n$$

$$\text{e.g. } D' = \{1, 2, 4, 6, 6\}$$

$$E[D'] = 3.8$$

## Mean Values Quiz

$$1. D = \{1, 2, 3\}, E[D] = \frac{1+2+3}{3} = 2$$

$$2. D = \left\{ \begin{bmatrix} 1 \\ 7 \end{bmatrix}, \begin{bmatrix} 2 \\ 8 \end{bmatrix}, \begin{bmatrix} 3 \\ 9 \end{bmatrix} \right\}, E[D] = \begin{bmatrix} \frac{1+2+3}{3} \\ \frac{4+8+6}{3} \\ \frac{7+8+9}{3} \end{bmatrix} = \begin{bmatrix} 2 \\ 5 \\ 8 \end{bmatrix}$$

$$3. D = \left\{ 2 \begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix}, 2 \begin{bmatrix} 3 \\ 4 \\ 5 \end{bmatrix}, 2 \begin{bmatrix} 5 \\ 1 \\ 3 \end{bmatrix} \right\}, E[D] = \begin{bmatrix} \frac{2+6+10}{3} \\ \frac{4+8+6}{3} \\ \frac{6+10+2}{3} \end{bmatrix} = \begin{bmatrix} 6 \\ 6 \\ 6 \end{bmatrix}$$

$$4. D = \left\{ \begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix} + \begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix}, \begin{bmatrix} 3 \\ 4 \\ 5 \end{bmatrix} + \begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix}, \begin{bmatrix} 5 \\ 3 \\ 1 \end{bmatrix} + \begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix} \right\}, E[D] = \begin{bmatrix} \frac{2+4+6}{3} \\ \frac{4+6+5}{3} \\ \frac{6+8+4}{3} \end{bmatrix} = \begin{bmatrix} 4 \\ 5 \\ 6 \end{bmatrix}$$

5.  $\bar{x}_{n-1}$  of dataset  $D_{n-1}$  with  $n-1$  data points

collect new data point  $x_n$

new mean  $\bar{x}_n$  of full dataset  $D_n = D_{n-1} \cup \{x_n\}$

$$\hookrightarrow \bar{x}_n = \bar{x}_{n-1} + \frac{1}{n}(x_n - \bar{x}_{n-1})$$

e.g.

$$D_{n-1} = \{1, 2, 3\}, x_n = 4, D_n = \{1, 2, 3, 4\}$$

$$E[D_{n-1}] = 2, E[D_n] = 2.5$$

$$2.5 = 2 + \frac{1}{4}(4-2)$$

6. 2D array of 28x28 image: reshape to vector of length 784 (`x.flatten()`)

## Variance of 1D datasets



$E[D_1] = E[D_2] = 3$  (same mean, but data points in  $D_2$  less concentrated around mean compared to  $D_1$ )

Variance = used to characterise variability / spread of points in data

$$D_1 = \{1, 2, 4, 5\}, E[D_1] = 3$$

$$D_2 = \{-1, 3, 7\}, E[D_2] = 3$$

↙ Sum of Sq. distances

$$D_1: \frac{(1-3)^2 + (2-3)^2 + (4-3)^2 + (5-3)^2}{4} = 2.5 \quad (\text{avg. Squared distance from mean value})$$

$$D_2: \frac{(-1-3)^2 + (3-3)^2 + (7-3)^2}{3} = \frac{32}{3}$$

↗ larger

(Spread is higher in  $D_2$ )

$$\{x_1, \dots, x_n\} =: X$$

$$\text{Var}[X] = \frac{1}{N} \sum_{n=1}^N (x_n - \mu)^2$$

- can never be -ve!
- expressed in squared units

$$\mu = E[X]$$

$$\text{Standard deviation} = \sqrt{\text{Var}[x]}$$

↑ same units as mean value

(usually use st. deviations when talking about spread of data)

## Variance of 1D datasets quiz

1.  $D = \{1, 2, 3, 2\}$ ,  $\text{Var}[D] = \frac{(1-2)^2 + (2-2)^2 + (3-2)^2 + (2-2)^2}{4} = 0.5$   
 $E[D] = 2$

2. Standard deviation =  $\sqrt{0.5} =$

3. Add 1 to each element in D :  $\text{Var}[D] = \frac{(2-3)^2 + (3-3)^2 + (4-3)^2 + (3-3)^2}{4} = 0.5$  (adding a constant doesn't change the variance)  
 $E[D] = 3$

4. Multiply each sample in D by 2 :  $D = \{2, 4, 6, 4\}$ ,  $\text{Var}[D] = \frac{(2-4)^2 + (4-4)^2 + (6-4)^2 + (4-4)^2}{4} = 2$ ,  $\therefore$  4 times variance of D  
 $E[D] = 4$

↙ variance

5.  $\bar{x}_{n-1}$ ,  $\sigma_{n-1}^2$  for  $D_{n-1}$  with  $n-1$  samples

What is variance  $\sigma_n^2$  if we add new element  $x_*$ :

$$\sigma_n^2 = \frac{n-1}{n} \sigma_{n-1}^2 + \frac{1}{n} (x_* - \bar{x}_{n-1})(x_* - \bar{x}_n)$$

e.g.  $D_{n-1} = \{1, 3, 5\}$ ,  $E[D_{n-1}] = 3$ ,  $\sigma_{n-1}^2 = \frac{(1-3)^2 + (3-3)^2 + (5-3)^2}{3} = 2.6$

$$x_* = 4 : D_n = \{1, 3, 5, 4\}, E[D_n] = \frac{13}{4}, \sigma_n^2 = \frac{(1-\frac{13}{4})^2 + (3-\frac{13}{4})^2 + (5-\frac{13}{4})^2 + (4-\frac{13}{4})^2}{4} = 2.1875$$

$$2.1875 = \frac{3}{4}(2.6) + \frac{1}{4}(4-3)(4-\frac{13}{4})$$

Symmetric, positive definite matrices

Matrix  $M$  is symmetric if  $M = M^T$

Symmetric matrix  $M$  is positive-definite if  $z^T M z > 0$ :

$$\begin{bmatrix} z_1 & z_2 & \dots & z_n \end{bmatrix} \begin{bmatrix} M_{11} & M_{12} & \dots & M_{1n} \\ M_{21} & M_{22} & \ddots & \vdots \\ \vdots & \ddots & \ddots & \vdots \\ M_{n1} & \dots & \dots & M_{nn} \end{bmatrix} \begin{bmatrix} z_1 \\ z_2 \\ \vdots \\ z_n \end{bmatrix} > 0$$

for a real column vector  $z \neq 0$

Symmetric and  
positive-definite  
 $\downarrow$   
(e.g.  $\begin{bmatrix} 7 & 1 \\ 1 & 2 \end{bmatrix}$ )

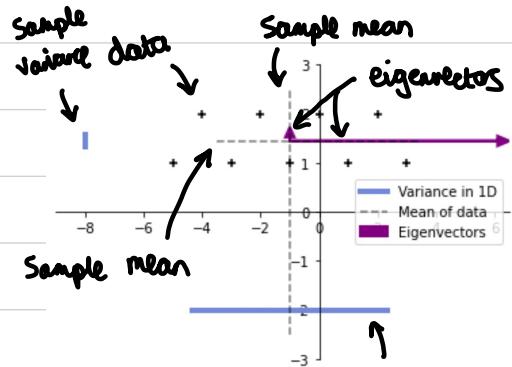
The eigenvalues of Symmetric positive-definite matrices are +ve

↳ e.g.  $M$  is Symmetric positive-definite matrix,  $\lambda$  is eigenvalue of  $M$ , so  $Mv = \lambda v$ , where  $v$  is the corresponding eigenvector

↳ when  $M$  is positive-definite, we have  $v^T M v = \lambda v^T v > 0$

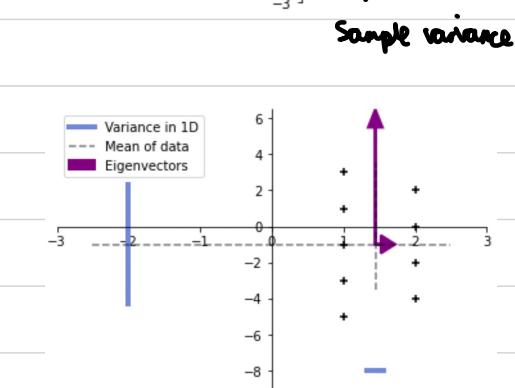
↳ since  $v^T v > 0$ ,  $\lambda$  must be +ve

Covariance matrix: Symmetric and positive-definite



Covariance matrix for dataset of 2 indep. r.v.  $x$  and  $y$

$$\begin{bmatrix} \text{cov}(x,x) & \text{cov}(x,y) \\ \text{cov}(y,x) & \text{cov}(y,y) \end{bmatrix} = \begin{bmatrix} 7.5 & 0 \\ 0 & 1.1 \end{bmatrix}$$



$$\begin{bmatrix} \text{cov}(x,x) & \text{cov}(x,y) \\ \text{cov}(y,x) & \text{cov}(y,y) \end{bmatrix} = \begin{bmatrix} 1.11 & 0 \\ 0 & 7.5 \end{bmatrix}$$

In both covariance matrices, 2 variables are independent since  $\text{cov}(x,y) = \text{cov}(y,x) = 0$  ( $\therefore$  matrices are diagonal)

$\hookrightarrow$  diagonal entries: variances for each of the variables, always true

$\hookrightarrow$  eigenvalues of matrix = variances

Eigenvectors are orthogonal to one another due to matrix symmetry

$\hookrightarrow$  represent different directions in which the data varies

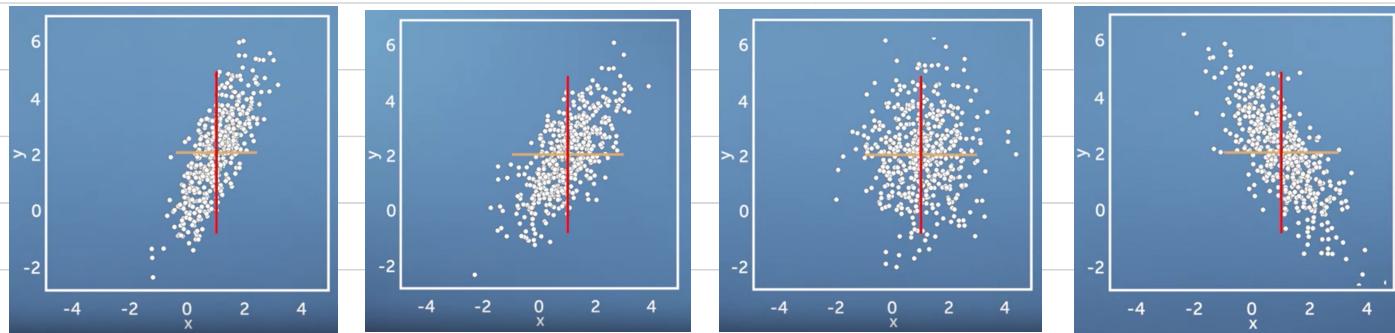
$$\hookrightarrow \text{cov}(x,y) = \text{cov}(y,x)$$

$\hookrightarrow$  directions in which variables vary the most coincide with the eigenvectors (in examples: along x and y axes)

Matrix of eigenvalues  $D$ : obtained by diagonalisation  $\Rightarrow D = C^{-1}MC$ ,  $C$ : matrix of eigenvectors

$\hookrightarrow$  covariance matrix similar to diagonal matrix with +ve, real entries (variances along the eigenvectors)

## Variance of higher dimensional datasets



Same variances in  $x$  and  $y$ , same

means, but different shapes

$\hookrightarrow$  if we look at horizontal and vertical spread of data: can't explain any correlation between  $x$  and  $y$

$$\text{Cov}[x, y] = E[(x - \mu_x)(y - \mu_y)]$$

$$\mu_x = E[x]$$

$$\mu_y = E[y]$$

$$\text{Var}[x]$$

$$\text{Var}[y]$$

$$\text{Cov}[x, y]$$

$$\text{Cov}[y, x]$$

$$\begin{bmatrix} \text{Var}[x] & \text{Cov}[x, y] \\ \text{Cov}[y, x] & \text{Var}[y] \end{bmatrix}$$

(covariance matrix)

always symmetric, positive definite matrix

Cov. between  $x$  and  $y$  is +ve: on avg.  $y$  increases if we increase  $x$

Cov. between  $x$  and  $y$  is -ve: on avg.  $y$  decreases if we increase  $x$

Cov. between  $x$  and  $y$  is 0:  $x$  and  $y$  are uncorrelated

Covariance matrix if every element is 3D vector with  $(x, y, z)$  components:

$$\begin{bmatrix} \text{Var}[x] & \text{Cov}[x,y] & \text{Cov}[x,z] \\ \text{Cov}[y,x] & \text{Var}[y] & \text{Cov}[y,z] \\ \text{Cov}[z,x] & \text{Cov}[z,y] & \text{Var}[z] \end{bmatrix} \quad (\text{variances on diagonal, cross-covariances on off-diagonal})$$

Dataset of  $D$ -dimensional vectors  $x = \begin{bmatrix} x_1 \\ \vdots \\ x_D \end{bmatrix}$ , covariance matrix:  $D \times D$

$$D = \{x_1, \dots, x_N\}, x_i \in \mathbb{R}^D$$

$$\text{Var}[D] = \frac{1}{N} \sum_{i=1}^N (x_i - \mu)(x_i - \mu)^T$$

$\uparrow$   
 $D \times D$  matrix

Covariance matrix of a 2D dataset quiz

$$1. \text{ 2D dataset } D = \{x_i\}_{i=1}^N \text{ with } N \text{ samples}$$

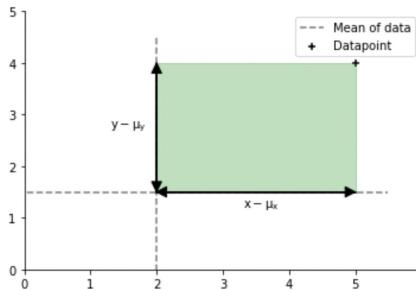
$\hookrightarrow$  each sample  $x_i$ : 2D vector of  $x, y$

$$\begin{bmatrix} 0 & 1 \\ 1 & 2 \\ 0 & 1 \\ 1 & 2 \end{bmatrix} \begin{bmatrix} x & y \\ y & x \end{bmatrix} \begin{bmatrix} -0.5 & -0.5 \\ 0.5 & 0.5 \\ -0.5 & -0.5 \\ 0.5 & 0.5 \end{bmatrix} \begin{bmatrix} x \\ y \end{bmatrix} \begin{bmatrix} -0.5 & 0.5 & -0.5 & 0.5 \\ -0.5 & 0.5 & -0.5 & 0.5 \end{bmatrix}$$

$(4 \times 2) \quad (2 \times 4)$   
 $\curvearrowleft$   
(other way around)

$$\text{Covariance between 2 scalar random variables: } \text{Cov}[x, y] = E[(x - \mu_x)(y - \mu_y)] \approx \frac{1}{N} \sum_{i=1}^N (x_i - \mu_x)(y_i - \mu_y)$$

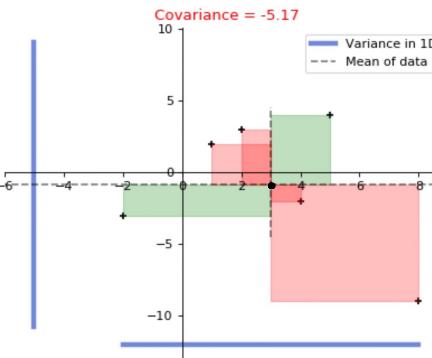
In the formula for covariance, we can think of each individual multiplication as the calculation of an area, a rectangle with sides  $x - \mu_x$  and  $y - \mu_y$ .



*as x increases from  $\mu_x$ ,  
y also increases*

For this datapoint, an increase in  $x$  from the mean is linked to an increase in  $y$ . Where  $x - \mu_x$  and  $y - \mu_y$  have the same sign, the contribution to the covariance is positive and in green, while if the signs are opposite it will be negative and in red. In other words, green means that  $x$  and  $y$  are positively correlated, while red means they're negatively correlated.

The total sum of areas, divided by the number of points  $n$ , will be the value of the covariance.

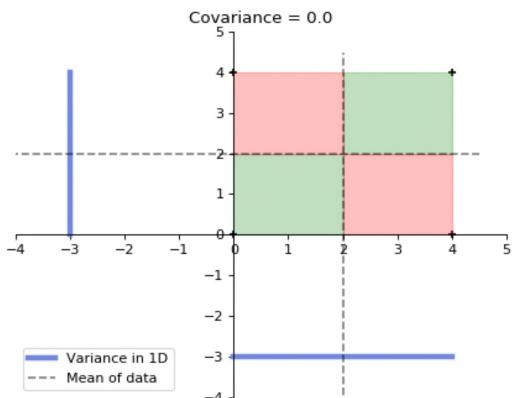


The dashed lines meet at the mean of the dataset. The blue lines represent the magnitude of the variance of the x (horizontal) and y (vertical) components of the dataset.

If red and green balance out, the covariance will be 0. Otherwise the sign of the covariance will give a direction in which the points appear to correlate.

What is  $\text{cov}(x, y)$  for the dataset in the array labelled "Q1: square"? Is it what you would expect from the plot?

*↳ covariance is 0*



*points are evenly distributed around the mean,  
So they balance out, so no way to determine  
direction of correlation*

$$2. \begin{bmatrix} \text{cov}(x,x) & \text{cov}(x,y) \\ \text{cov}(y,x) & \text{cov}(y,y) \end{bmatrix} = \begin{bmatrix} \text{Var}(x) & \text{cov}(x,y) \\ \text{cov}(y,x) & \text{Var}(y) \end{bmatrix}$$

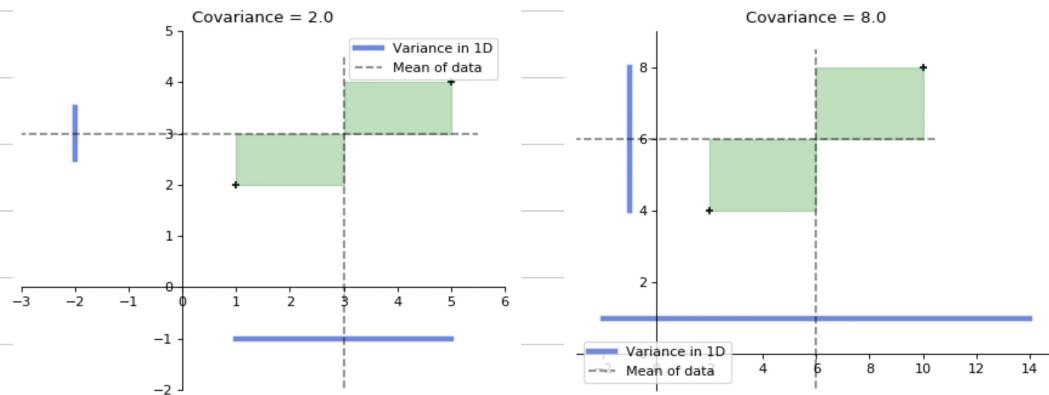
$$D = \left\{ \begin{bmatrix} 1 \\ 2 \end{bmatrix}, \begin{bmatrix} 5 \\ 4 \end{bmatrix} \right\}, \text{cov}[x,y] = \frac{1}{N} \sum_{i=1}^N (x - \mu_x)(y - \mu_y) = \frac{1}{2} [(1-3)(2-3) + (5-3)(4-3)] = 2$$

$$\mu_x = 3, \mu_y = 3$$

$$\text{Var}[x] = \frac{1}{2} [(1-3)^2 + (5-3)^2] = 4$$

$$\text{Var}[y] = \frac{1}{2} [(2-3)^2 + (4-3)^2] = 1$$

$$\text{Covariance matrix} = \begin{bmatrix} 4 & 2 \\ 2 & 1 \end{bmatrix}$$



$$3. \text{ Dataset with covariance matrix } \begin{bmatrix} 3 & 2 \\ 2 & 4 \end{bmatrix}$$

\* each element  
in D by 2

$$\begin{aligned} \text{Cov}[x,y] &\uparrow \times 4 \\ \text{Var}[x] &\uparrow \times 4 \\ \text{Var}[y] &\uparrow \times 4 \end{aligned}$$

$$\text{Multiply each vector in } D \text{ by 2: } \text{Var}[x'] = \text{Var}[x] \times 2^2 = 12$$

$$\text{Var}[y'] = \text{Var}[y] \times 2^2 = 16 \quad \Rightarrow \quad \begin{bmatrix} 12 & 8 \\ 8 & 16 \end{bmatrix}$$

$$\text{Cov}[x,y] = 2 \times 2^2 = 8$$

$$4. D = \left\{ \begin{bmatrix} 1 \\ 2 \end{bmatrix}, \begin{bmatrix} 7 \\ 4 \end{bmatrix} \right\}, \text{ covariance matrix } \begin{bmatrix} 9 & 3 \\ 3 & 1 \end{bmatrix}$$

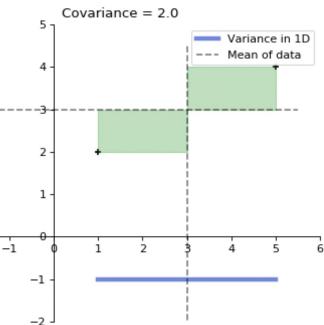
↪ add  $\begin{bmatrix} 2 \\ 2 \end{bmatrix}$  to each element in  $D$ :  $D = \left\{ \begin{bmatrix} 3 \\ 4 \end{bmatrix}, \begin{bmatrix} 9 \\ 6 \end{bmatrix} \right\}$

$$\text{Cov}[x, y] = \text{Cov}[y, x] = \frac{1}{2} [(3-6)(4-5) + (9-6)(6-5)] = 3$$

$$\mu_x = 6, \mu_y = 5$$

$$\text{Var}[x] = \frac{1}{2} [(3-6)^2 + (9-6)^2] = 9 \Rightarrow \begin{bmatrix} 9 & 3 \\ 3 & 1 \end{bmatrix}$$

$$\text{Var}[y] = \frac{1}{2} [(4-5)^2 + (6-5)^2] = 1$$



(covariance matrix  
doesn't change)

5. Dataset  $D$ , every element in  $D$  has  $x$  and  $y$  coordinate

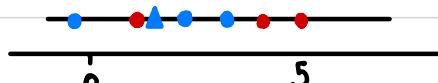
$$\hookrightarrow \text{covariance matrix} = \begin{bmatrix} 1 & 0.8 \\ 0.8 & 1 \end{bmatrix}$$

↪  $x$  and  $y$  are positively correlated - when  $x$  increases then  $y$  increases on average and vice versa

## Effect on the mean

Linear transformation: Shift data around or stretch it

$$D = \{-1, 2, 3\}, E[D] = \frac{4}{3}$$



Shift blue data to right by 2: mean also shifts by 2

$$D' = \{1, 4, 5\} = D + 2, E[D'] = \frac{10}{3} = \frac{4}{3} + 2$$

(mean + shift)

$$E[D+a] = a + E[D] \quad (\text{Shift by } a)$$

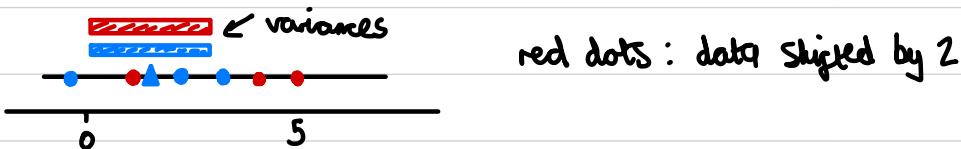
$$D'' = \{-2, 4, 6\}$$

$$E[D''] = \frac{8}{3} = \frac{4}{3} \cdot 2 \quad \begin{matrix} \leftarrow \text{Scaling factor} \\ \leftarrow \text{Scaling factor} \end{matrix}$$

$$E[\alpha D] = \alpha E[D] \quad (\text{Scale by } \alpha)$$

$$E[\alpha D + a] = \alpha E[D] + a \quad \begin{matrix} \leftarrow \text{Shift} \\ \uparrow \quad \leftarrow \text{Scaling factor} \end{matrix} \quad (\text{linear transformations on mean of data})$$

## Effect on the covariance



red dots : data shifted by 2

Shifting data: no effect on variance

$$\text{Var}[D] = \text{Var}[D+a]$$

$\uparrow$  offset applied to each element of  $D$

$$\text{Var}[\alpha D] = \alpha^2 \text{Var}[D] \quad (\text{e.g. Scale data by 2, squared distance from } \mu \text{ scaled by 4, so Var. } 4x \uparrow \text{ as before})$$

$\uparrow$  Scales each element, real number

$$D = \{x_1, \dots, x_N\}, x_i \in \mathbb{R}^p$$

Variance of data in  $D$ : given by covariance matrix

Apply linear transformation:  $Ax_i + b$

$$\text{Var}[AD + b] = A \text{Var}[D] A^T$$

$$\hookrightarrow \text{Var}[D] = E[(D - \mu)(D - \mu)^T]$$

$$\begin{aligned}\hookrightarrow \text{multiply } D \text{ by } A: \text{Var}[AD] &= E[(AD - A\mu)(AD - A\mu)^T] \\ &= E[A(D - \mu)(D - \mu)^T A^T] \\ &= A E[(D - \mu)(D - \mu)^T] A^T \\ &= A \text{Var}[D] A^T\end{aligned}$$

Linear transformation of data: Shifting only affects mean, scaling affects mean and variance

## Week 1 Lab

Means and (co)variances of a dataset

Vectorized code significantly improves runtime on larger datasets

e.g. data matrix  $X$  of size  $(N, D)$

↳ what is mean and covariance when we apply affine transformation  $Ax_i + b$  for each  $x_i$  in  $X$

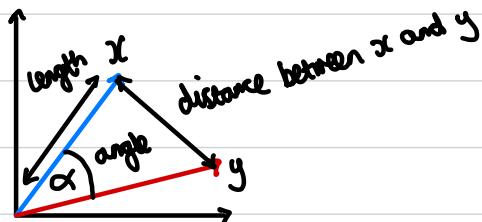
$$m' = \text{affine\_mean}(m, A, b) = A @ \text{mean} + b$$

$$S' = \text{affine\_covariance}(S, A, b) = A @ S @ A.T \quad (\text{covariance not affected by } b)$$

( $m$ : old mean,  $S$ : old covariance)

## Week 2: Inner Products

DR: find compact representations of data that live in a lower dimensional space (but similar to original data)

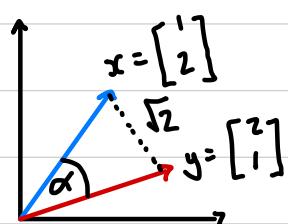


$$x^T y = \sum_{i=1}^N x_i y_i, \quad x, y \in \mathbb{R}^N$$

$$\|x\| = \sqrt{x^T x} = \sqrt{\sum_{i=1}^N x_i^2}$$

$$\text{e.g. } \|x\| = \sqrt{1^2 + 2^2} = \sqrt{5}, \quad \|y\| = \sqrt{5}$$

$$\cos \alpha = \frac{x^T y}{\|x\| \|y\|} = \frac{4}{\sqrt{5} \times \sqrt{5}} = \frac{4}{5}, \quad \therefore \alpha \approx 0.64 \text{ rad}$$



$$\begin{aligned} d(x-y) &= \left\| \begin{bmatrix} 1 \\ 2 \end{bmatrix} - \begin{bmatrix} 2 \\ 1 \end{bmatrix} \right\| \\ &= \left\| \begin{bmatrix} -1 \\ 1 \end{bmatrix} \right\| \\ &= \sqrt{1+1} \\ &= \sqrt{2} \end{aligned}$$

$$d(x, y) = \|x - y\| = \sqrt{(x-y)^T (x-y)}$$

Inner product: generalisation of the dot product

↪ want to express geometric properties (lengths, angles) between vectors

Def:  $x, y \in V$

$\langle \cdot, \cdot \rangle : V \times V \rightarrow \mathbb{R}$  (mapping from  $V \times V$  to real numbers)

• symmetric:  $\langle x, y \rangle = \langle y, x \rangle$

if and only if

• positive definite:  $\langle x, x \rangle \geq 0$ ,  $\langle x, x \rangle = 0 \Leftrightarrow x = 0$  ← zero vector

• bilinear:  $x, y, z \in V$ ,  $\alpha \in \mathbb{R}$ , inner product between  $\alpha x + z$  and  $y$ :  $\langle \alpha x + z, y \rangle = \alpha \underbrace{\langle x, y \rangle} + \langle z, y \rangle$

$$\langle x, \alpha y + z \rangle = \alpha \langle x, y \rangle + \langle x, z \rangle$$

Inner product:

•  $\langle x, y \rangle = x^T I y$  (dot product)

(bilinear: linearity in both arguments of function)

•  $\langle x, y \rangle = x^T A y$

$$\text{e.g. } A = \begin{bmatrix} 2 & 1 \\ 1 & 2 \end{bmatrix} \rightarrow 2x_1y_1 + x_1y_2 + x_2y_1 + 2x_2y_2$$

any symmetric positive-definite matrix defines a valid inner product

Lengths and distances

$$\|x\| = \sqrt{\langle x, x \rangle} \quad \text{inner product is positive-definite (i.e. } \geq 0, \text{ so can sqrt)}$$

↑ length / norm of  $x$

e.g.  $x = \begin{bmatrix} 1 \\ 1 \end{bmatrix}$

$$\langle x, y \rangle = x^T y \Rightarrow \|x\| = \sqrt{2} \quad (\text{v1})$$

or  $\langle x, y \rangle = x^T \begin{bmatrix} 1 & -0.5 \\ -0.5 & 1 \end{bmatrix} y = x_1 y_1 - \frac{1}{2}(x_1 y_2 + x_2 y_1) + x_2 y_2 \quad (\text{v2})$

$$\Rightarrow \|x\| = \sqrt{x_1^2 - \frac{1}{2}(x_1 y_2 + x_2 y_1) + x_2^2} = \sqrt{x_1^2 - x_1 x_2 + x_2^2}$$

$$\|x\|^2 = \langle x, x \rangle = 1+1-1=1, \Rightarrow \|x\|=1 \quad (\text{length of vector using v2 of inner product})$$

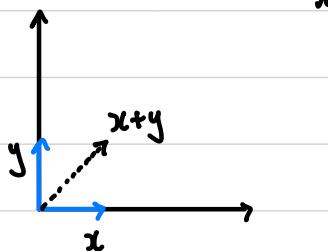
(Squared norm: inner product

of  $x$  with itself)

Stretch vector by  $\lambda$

$$\|\lambda x\| = |\lambda| \|x\| \quad \text{absolute}$$

Triangle inequality:  $\|x+y\| \leq \|x\| + \|y\|$



$$x = \begin{bmatrix} 1 \\ 0 \end{bmatrix}, y = \begin{bmatrix} 0 \\ 1 \end{bmatrix}$$

$$\|x\| = 1 = \|y\|$$

$$\|x+y\| = \sqrt{2}$$

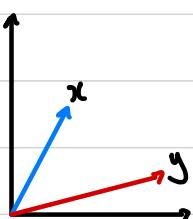
$$\sqrt{2} \leq 2$$

$$d(x, y) = \|x-y\| = \sqrt{\langle x-y, x-y \rangle} \quad (\text{length of difference vector})$$

use a dot product: distance = Euclidean distance

e.g.  $x = \begin{bmatrix} 2 \\ 3 \end{bmatrix}, y = \begin{bmatrix} 4 \\ 1 \end{bmatrix}$

$$x-y = \begin{bmatrix} -2 \\ 2 \end{bmatrix}$$



Dot product:  $\sqrt{8}$

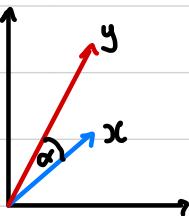
or:  $\langle x, y \rangle = x^T \begin{bmatrix} 1 & -0.5 \\ -0.5 & 1 \end{bmatrix} y \quad (\text{different inner product})$

$$\|x-y\| = \sqrt{12}$$

## Angles and Orthogonality

$$\cos \alpha = \frac{\langle x, y \rangle}{\|x\| \|y\|}, \text{ e.g. } x = \begin{bmatrix} 1 \\ 1 \end{bmatrix}, y = \begin{bmatrix} 1 \\ 2 \end{bmatrix}$$

$$\cos \alpha = \frac{x^T y}{\sqrt{x^T x} \sqrt{y^T y}} = \frac{3}{\sqrt{10}}$$

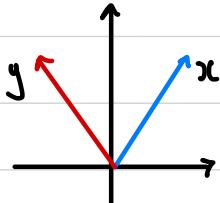


$$\therefore \alpha \approx 0.32 \text{ rad} \approx 18^\circ$$

$$\text{e.g. } x = \begin{bmatrix} 1 \\ 1 \end{bmatrix}, y = \begin{bmatrix} -1 \\ 1 \end{bmatrix}$$

$$\cos \alpha = 0, \therefore \alpha = \frac{\pi}{2} \text{ rad} = 90^\circ$$

(orthogonal vectors)



Inner product allows us to characterise orthogonality : vectors which are most dissimilar, nothing in common besides origin

2 vectors  $x$  and  $y$ , where  $x$  and  $y$  are non-zero vectors are orthogonal iff their inner product = 0

Vectors orthogonal w.r.t. 1 inner product don't have to be orthogonal w.r.t another inner product

e.g. different inner product:  $\langle x, y \rangle = x^T \begin{bmatrix} 2 & 0 \\ 0 & 1 \end{bmatrix} y = 7 \quad \langle x, y \rangle = -1 \quad (\text{vectors not orthogonal w.r.t. this inner product})$

Can find basis of vector space so basis are orthogonal to each other

↳  $\langle b_i, b_j \rangle = 0$ , if  $i \neq j$

↳ use inner product to normalize basis vectors:  $\|b_i\| = 1$  (orthonormal basis)

## Unconventional Inner Products

Inner product between 2 functions:  $\langle u, v \rangle = \int_a^b u(x)v(x) dx$  ( $\text{if } \int f^2 dx = 0, \text{ functions } u \text{ and } v \text{ are orthogonal}$ )

$$\text{e.g. } u(x) = \sin(x)$$

$$v(x) = \cos(x)$$

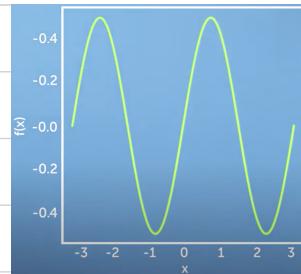
$$f(x) = u(x)v(x) \rightarrow$$

if we set integral limit to

$-\pi$  and  $\pi$ , integral of

$$f(x) = 0, \text{ so } \sin \text{ and } \cos$$

are orthogonal



odd function, as  $f(-x) = -f(x)$

e.g. Set of functions  $\{1, \cos x, \cos 2x, \cos 3x, \dots\}$

are all orthogonal if we integrate from  $-\pi$  to  $\pi$

2 uncorrelated r.v.:  $\text{Var}[x+y] = \underbrace{\text{Var}[x]}_{\text{r.v.}} + \underbrace{\text{Var}[y]}_{\text{r.v.}}$

Pythagorean theorem:  $c^2 = a^2 + b^2$

r.v. can be considered vectors in a vector space

↙ Symmetric, positive-definite, linear

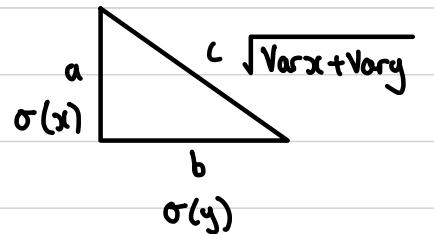
Inner product between r.v. :  $\langle x, y \rangle = \text{Cov}[x, y]$

$$\text{Cov}[\alpha x + y, z] = \alpha \text{Cov}[x, z] + \text{Cov}[y, z] \quad (\text{linearity})$$

$$\|x\| = \text{length of r.v.} : \sqrt{\text{Cov}[x, x]} = \sqrt{\text{Var}[x]} = \sigma(x)$$

↪ 0 vector has no uncertainty, as  $\sigma$  is 0

Angle between 2 r.v. :  $\cos \theta = \frac{\langle x, y \rangle}{\|x\| \|y\|} = \frac{\text{Cov}[x, y]}{\sqrt{\text{Var}[x] \text{Var}[y]}}$  ( $= 0$  iff  $\text{Cov}[x, y] = 0$ , is the case when  $x$  and  $y$  are uncorrelated)



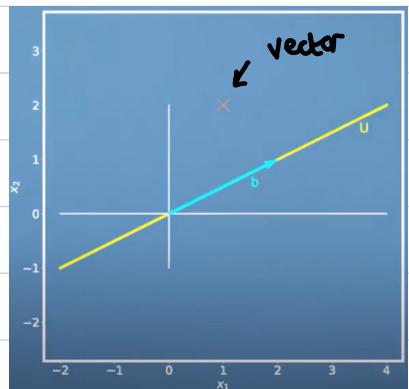
(geometric interpretation of r.v.)

## Week 3 : Orthogonal Projections

Compress data : loss of information

↳ want to keep informative dimensions, ignore irrelevant dimensions

### Projection onto 1D Subspaces

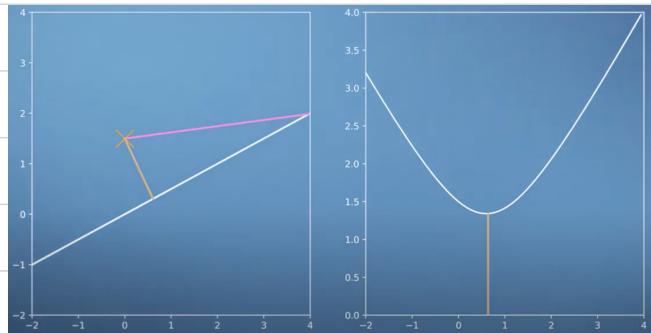


Vector can be represented as linear combination of basis vectors

1D Subspace  $U$  with basis vector  $b$

↳ all vectors in  $U$  : represented as  $\lambda b$  for some lambda

Want a vector in  $U$  that is closest to  $x$



find vector in  $U$  closest to  $x$  by an orthogonal projection of  $x$  onto  $U$  (difference vector of  $x$  and projection is orthogonal to  $U$ )

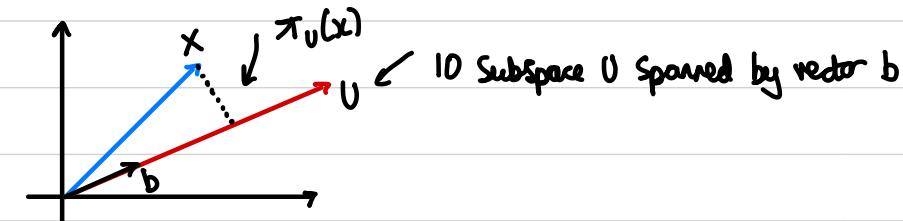
↙ projection of  $x$  onto  $U$

$\pi_U(x)$

1.  $\pi_U(x) \in U \Rightarrow \exists s \in \mathbb{R} : \pi_U(x) = sb$  (as  $\pi_U(x) \in U$ )

2.  $\langle b, \pi_U(x) - x \rangle = 0$  (orthogonality condition: inner product of basis vector and difference vector is 0)

(both properties generally hold for any  $x$  in  $\mathbb{R}^D$  and 1D Subspaces  $U$ )



$$\langle b, \pi_U(x) - x \rangle = 0$$

$$\Leftrightarrow \langle b, \overline{\pi_U(x)} \rangle - \langle b, x \rangle = 0 \quad (\text{exploit linearity of inner product})$$

$$\Leftrightarrow \langle b, s b \rangle - \langle b, x \rangle = 0$$

$$\Leftrightarrow s \|b\|^2 - \langle b, x \rangle = 0$$

$$\Leftrightarrow s = \frac{\langle b, x \rangle}{\|b\|^2}$$

$$\Rightarrow \pi_U(x) = sb = \frac{\langle b, x \rangle b}{\|b\|^2}$$

↳ choose dot product as inner product:  $\frac{\overbrace{\langle b^T, x \rangle b}^{\text{dot product}}}{\|b\|^2} = \underbrace{\frac{b b^T}{\|b\|^2} x}_{\text{projection matrix}} = \pi_U(x)$

projection matrix: projects 2D point onto 1D Subspace

( $s$ : coordinate of projection w.r.t.  $b$  of Subspace  $U$ )

↙ multiple of basis vector that spans  $U$

coordinate of projected point w.r.t. basis  $b$  by using dot product  
of  $b$  with  $x$

$$\text{if } \|b\|=1 \Rightarrow s = b^T x, \pi_U(x) = \underline{b b^T x}$$

projection matrix

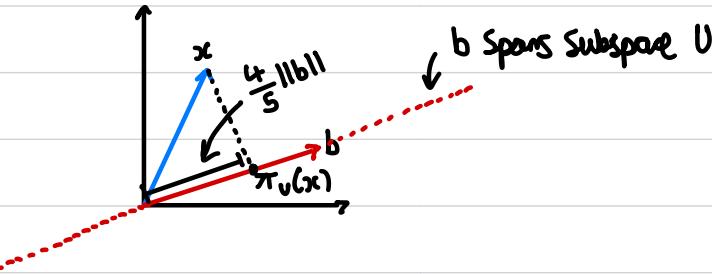
Projection  $\pi_U(x)$  still a vector in  $\mathbb{R}^D$

↳ but don't need  $D$  coordinates to represent it, only  
need  $s$

## Example 1D Projection

$$\pi_U(x) = \frac{x^T b}{\|b\|^2} b$$

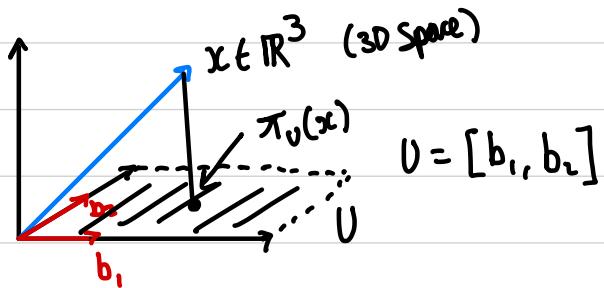
e.g.  $b = \begin{bmatrix} 2 \\ 1 \end{bmatrix}$ ,  $x = \begin{bmatrix} 1 \\ 2 \end{bmatrix}$



want orthogonal projection of  $x$  onto  $U$

$$\pi_U(x) = \frac{4}{5} \begin{bmatrix} 2 \\ 1 \end{bmatrix} \text{ (orthogonal projection)}$$

## N-dimensional Projections



$\pi_U(x)$  is element of  $U$ , so can represent as linear combination of basis vectors of  $U$

$$\pi_U(x) = \gamma_1 b_1 + \gamma_2 b_2$$

$$\langle x - \pi_U(x), b_1 \rangle = 0$$

$$\langle x - \pi_U(x), b_2 \rangle = 0$$

represent projection as L.C. of basis of Subspace

$$1. \pi_U(x) = \sum_{i=1}^M \gamma_i b_i$$

$$\boldsymbol{\gamma} = \begin{bmatrix} \gamma_1 \\ \vdots \\ \gamma_M \end{bmatrix}, \quad B = [b_1 | \dots | b_M]$$

$$2. \langle \pi_U(x) - x, b_i \rangle = 0, \quad i=1, \dots, M$$

$$\pi_U(x) = B\boldsymbol{\gamma}$$

difference vector between  $x$  and projection is orthogonal to Subspace

$$\downarrow \text{use dot product as inner product: } \langle \pi_U(x) - x, b_i \rangle = \langle B\boldsymbol{\gamma} - x, b_i \rangle = 0 \quad ? \quad \leftarrow \text{needs to be 0}$$

$$\Leftrightarrow \langle B\boldsymbol{\gamma}, b_i \rangle - \langle x, b_i \rangle = 0, \quad i=1, \dots, M$$

$$\Leftrightarrow \boldsymbol{\gamma}^T B^T b_i - x^T b_i = 0, \quad i=1, \dots, M$$

$$\Leftrightarrow \boldsymbol{\gamma}^T B^T B - x^T B = 0$$

$$\Leftrightarrow \boldsymbol{\gamma}^T = x^T B (B^T B)^{-1}$$

$$\Leftrightarrow \boldsymbol{\gamma} = (B^T B)^{-1} B^T x$$

$$\Rightarrow \pi_U(x) = B\boldsymbol{\gamma} = B \underbrace{(B^T B)^{-1} B^T}_{\text{projection matrix}} x$$

projection matrix

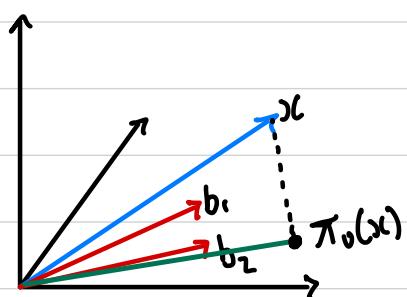
Special case of orthonormal basis:  $B^T B = I$

$$\Rightarrow \pi_U(x) = B B^T x$$

↑ projected vector in  $D$ , but only require  $M$  coordinates ( $\boldsymbol{\gamma}$  vector) to be represented as linear combination of basis vectors of Subspace  $U$

$$\text{1D case: } \pi_U(x) = \frac{b^T x}{b^T b} b, \quad \boldsymbol{\gamma} = \frac{b^T x}{b^T b}$$

## Example N-dimensional Projections



$$x = \begin{bmatrix} 2 \\ 1 \\ 1 \end{bmatrix}, b_1 = \begin{bmatrix} 1 \\ 2 \\ 0 \end{bmatrix}, b_2 = \begin{bmatrix} 1 \\ 0 \\ 1 \end{bmatrix}$$

$$U = [b_1, b_2]$$

Projected vector can be represented as L.C. of basis of Subspace, and vector that connects data point and its projection must be orthogonal to the subspace

Orthogonal projection :  $\pi_U(x) = Bz$

$$B = [b_1 \mid b_2] = \begin{bmatrix} 1 & 1 \\ 2 & 0 \\ 0 & 1 \end{bmatrix}$$

$$z = (B^T B)^{-1} B^T x$$

$$B^T x = \begin{bmatrix} 4 \\ 3 \\ 3 \end{bmatrix}$$

$$B^T B = \begin{bmatrix} 5 & 3 \\ 3 & 2 \end{bmatrix}$$

$$B^T B z = B^T x$$

↙ projected point has 3<sup>rd</sup> component of 0, Subspace requires 3<sup>rd</sup> component is always 0

$$(\text{Gaussian elimination}) \Rightarrow z = \begin{bmatrix} -1 \\ 3 \end{bmatrix}$$

$$\Rightarrow \pi_U(x) = -1b_1 + 3b_2 = \begin{bmatrix} 2 \\ 1 \\ 0 \end{bmatrix}$$

↙ 3D vector, but can represent using 2 coordinates if we use basis defined by  $b_1$  and  $b_2$ ,  $\therefore$  is compact representation of projection of  $x$  onto lower dimensional Subspace

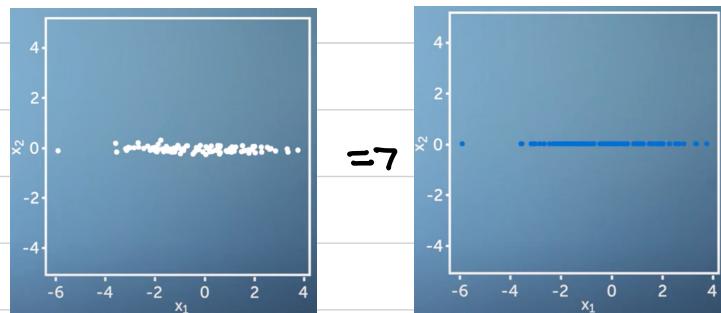
## Week 4 : Principal Component Analysis (PCA)

Algorithm for linear dimensionality reduction

↳ used for data compression + visualisation

High dimensional data (e.g. images) often lie in lower dimensional Subspace

↳ many dimensions are highly correlated



PCA: use orthogonal projections to find lower dimensional representations of data that retain as much info as possible expressed using fewer basis

data  $\rightarrow X = \{x_1, \dots, x_N\}$ ,  $x_i \in \mathbb{R}^D$   $\leftarrow$  D-dimensional vectors  
vector

use dot product for inner product

usual orthonormal basis of  $\mathbb{R}^D$

$$\textcircled{1} \quad x_n = \sum_{i=1}^D \beta_{in} b_i \quad (\text{every vector in } \mathbb{R}^D \text{ can be represented as L.C. of basis vectors})$$

$$\textcircled{2} \quad \beta_{in} = x_n^T b$$

orthogonal projection of  $x_n$  onto 1D Subspace Spanned by  $i^{\text{th}}$  basis vector

$$\textcircled{3} \quad B = (b_1, \dots, b_M) \quad (\text{matrix of orthonormal basis vectors})$$

$$\tilde{x} = \underline{BB^T x} \quad (\text{orthogonal projection of } x \text{ onto Subspace Spanned by } M \text{ basis vectors})$$



coordinates of  $\tilde{x}$  w.r.t. basis vectors in  $B$  ("code")

1. Centred data:  $E[X] = 0$

2. ONB  $b_1, \dots, b_D$

in PCA: ignore this

$$\tilde{x}_n = \sum_{i=1}^M \beta_{in} b_i + \sum_{i=M+1}^D \beta_{in} b_i \in \mathbb{R}^D$$

$M$ -dimensional Subspace ( $D-M$ )-dimensional Subspace (orthogonal complement to first subspace)

(2 sums from property ①)



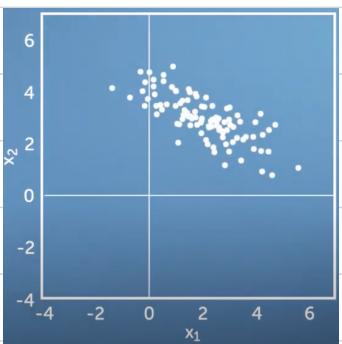
: principal Subspace ( $b_1, \dots, b_M$  Span this)

$\tilde{x}_n$  still D-dimensional vector, but lies in M-dimensional subspace of  $\mathbb{R}^D$

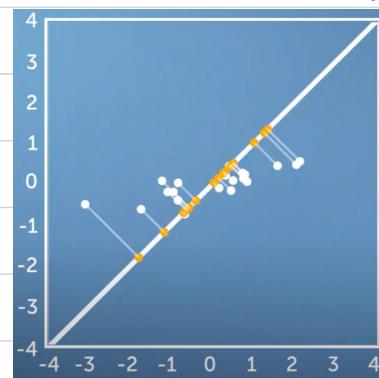
↳ only M coordinates ( $\beta_{n1}, \dots, \beta_{nM}$ ) necessary to represent it

Want to find parameters  $\beta_{in}$ , orthonormal basis vectors  $b_i$  such that avg. Squared reconstruction error is minimised

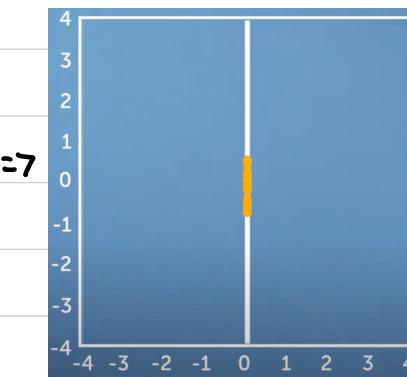
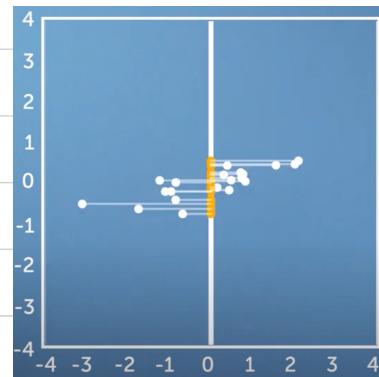
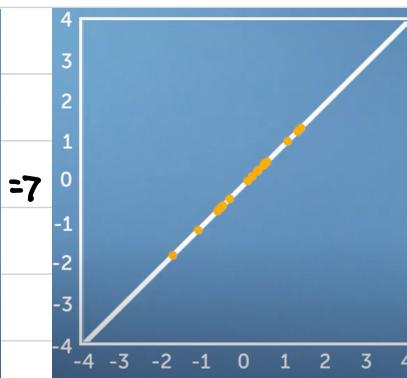
$$J = \frac{1}{N} \sum_{n=1}^N \|x_n - \tilde{x}_n\|^2 \quad (\text{loss function})$$



want to find good  
1D Subspace such  
that avg. Squared  
reconstruction error  
of data and projections  
is minimised



(options of subspaces)



Some projections more informative  
than others : PCA finds the  
best one

↳ partial deriv. of  $J$  w.r.t.  
parameters  $\beta_{in}$  and  $b_i$ , Set  
to 0, Solve for optimal  
parameters

$$\frac{\partial J}{\partial \{B_{in}, b_i\}} = \frac{\partial J}{\partial \tilde{x}_n} \frac{\partial \tilde{x}_n}{\partial \{B_{in}, b_i\}}$$

$$\frac{\partial J}{\partial \tilde{x}_n} = -\frac{2}{N} (x_n - \tilde{x}_n)^T$$

$$\frac{\partial J}{\partial B_{in}} = \frac{\partial J}{\partial \tilde{x}_n} \frac{\partial \tilde{x}_n}{\partial B_{in}}, \quad \frac{\partial \tilde{x}_n}{\partial B_{in}} = b_i, \quad i=1 \dots M$$

optimal coordinates of  $\tilde{x}_n$  wr.t.

basis are orthogonal projections of original data

point onto  $i$ 'th basis vector that spans principal subspace

$$\begin{aligned} & \downarrow \\ &= -\frac{2}{N} (x_n - \tilde{x}_n)^T b_i \\ &= -\frac{2}{N} (x_n - \sum_{j=1}^M B_{in} b_j)^T b_i \stackrel{\text{ONB}}{=} -\frac{2}{N} (x_n^T b_i - B_{in} \underbrace{b_i^T b_i}_{=1}) = -\frac{2}{N} (x_n^T b_i - B_{in}) = 0 \Leftrightarrow B_{in} = x_n^T b_i \end{aligned}$$

Ⓐ  $\tilde{x}_n = \sum_{j=1}^M B_{in} b_j$  (projected data point)

Ⓑ  $J = \frac{1}{N} \sum_{n=1}^N \|x_n - \tilde{x}_n\|^2$  (loss function)

Ⓒ  $\frac{\partial J}{\partial \tilde{x}_n} = -\frac{2}{N} (x_n - \tilde{x}_n)^T$  (partial deriv. of loss wr.t.  $\tilde{x}_n$ )

Ⓓ  $B_{in} = x_n^T b_i, \quad i=1, \dots, M$  (optimal coordinates)

$$\textcircled{A} \quad \tilde{x}_n = \sum_{j=1}^M \beta_{jn} b_j$$

$$\textcircled{B} \quad = \sum_{j=1}^M (\tilde{x}_n^T b_j) b_j \quad (\text{dot product is symmetric})$$

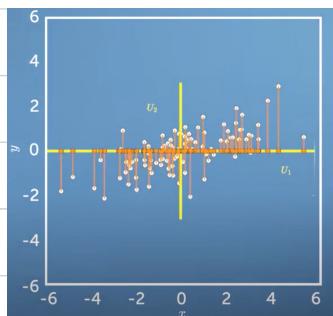
$$= \sum_{j=1}^M b_j (\underbrace{b_j^T \tilde{x}_n}_{\text{projection matrix}}) = \underbrace{\left( \sum_{j=1}^M b_j b_j^T \right)}_{Q} \tilde{x}_n \quad Q \quad (\tilde{x}_n: \text{orthogonal projection of } x_n \text{ onto Subspace Spanned by } M \text{ basis vectors } b_j \text{ from } j=1 \dots M)$$

$$x_n = \left( \sum_{j=1}^M b_j b_j^T \right) x_n \quad (\text{projection onto principal Subspace})$$

$$+ \left( \sum_{j=M+1}^D b_j b_j^T \right) x_n \quad (\text{projection onto orthogonal complement: missing from } \text{, } \therefore \text{ approximation})$$

difference between  
 $x_n$  and projection

$$\begin{aligned} x_n - \tilde{x}_n &= \left( \sum_{j=N+1}^D b_j b_j^T \right) x_n \quad (\text{displacement vector lies exclusively in subspace we ignore: orthogonal complement to} \\ &\quad \text{principal Subspace}) \\ &= \sum_{j=N+1}^D (b_j^T x_n) b_j \quad \textcircled{E} \end{aligned}$$



project data onto  $U_1$  Subspace

vertical lines: difference between data and projection (difference vector)

$\hookrightarrow$  have no variation in  $x$ , only has component that lies in  $U_2$  Subspace which is orthogonal complement to  $U_1$ , which is subspace we projected onto

$$\text{Regenerate loss function: } J = \frac{1}{N} \sum_{n=1}^N \|x_n - \tilde{x}_n\|^2 \quad \swarrow \text{avg. Squared reconstruction error}$$

$$\stackrel{(E)}{=} \frac{1}{N} \sum_{n=1}^N \left\| \sum_{j=M+1}^D (b_j^T x_n) b_j \right\|^2$$

$$\stackrel{\text{ONB}}{=} \frac{1}{N} \sum_{n=1}^N \sum_{j=M+1}^D (b_j^T x_n)^2$$

$$= \frac{1}{N} \sum_n \sum_j b_j^T x_n x_n^T b_j$$

$$= \sum_{j=M+1}^D b_j^T \left( \frac{1}{N} \sum_{n=1}^N x_n x_n^T \right) b_j$$

$\underbrace{\quad}_{\text{data covariance}}$

matrix  $S$  (we assume  
centered data)

$$J = \sum_{j=M+1}^D b_j^T S b_j = \text{trace} \left( \left( \sum_{j=M+1}^D b_j b_j^T \right) S \right)$$

(F)

$\underbrace{\quad}_{\text{projection matrix}}$

$\uparrow$  projects data covariance matrix onto orthogonal complement of principal Subspace

$\downarrow$  can regenerate loss function as variance of data projected onto Subspace we ignore

$\downarrow \therefore$  minimising (B) = minimising variance of data that lies in Subspace orthogonal to principal Subspace

$\downarrow$  interested in retaining as much variance after projection as possible

Minimising avg. Squared reconstruction error = minimising projection of variance of data by projecting onto Subspace we will ignore in PCA

$$J = \sum_{j=M+1}^D b_j^T S b_j \quad (\text{minimising this: need orthonormal basis that spans Subspace we will ignore})$$

↳ then take orthogonal complement as basis of principal Subspace

$$b_1, b_2 \quad b_i^T b_j = \delta_{ij} \quad (1 \text{ if } i=j, 0 \text{ otherwise})$$

$$J = b_2^T S b_2, \quad b_2^T b_2 = 1 \quad (\text{constraint})$$

Solve optimisation using Lagrangian:  $L = b_2^T S b_2 + \lambda (1 - b_2^T b_2)$

$$\frac{\partial L}{\partial \lambda} = 1 - b_2^T b_2 = 0 \Leftrightarrow b_2^T b_2 = 1$$

$$\frac{\partial L}{\partial b_2} = 2b_2^T S - 2\lambda b_2^T = 0 \Leftrightarrow Sb_2 = \lambda b_2 \quad (\text{eigenvalue problem})$$

$\lambda$  eigenvalue

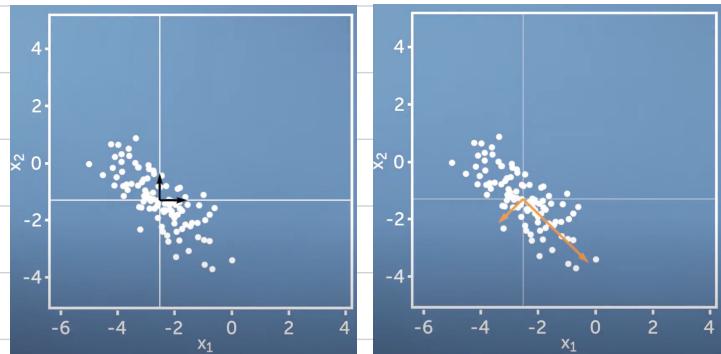
need to choose  $b_2$  as corresponding eigenvector which  
will span Subspace we will ignore

$b_2$  spans principal Subspace: eigenvector  
corresponding to largest  $\lambda$  of covariance matrix

$$J = b_2^T S b_2 = \underbrace{b_2^T b_2}_{\text{orthonormal basis}} \lambda = \lambda$$

∴ avg. Squared reconstruction error is minimised if  $\lambda$  is smallest  
eigenvalue of data covariance matrix

eigenvectors of covariance matrix already orthogonal due to matrix symmetry



best projection we can get which retains most info.

projects onto subspace spanned by eigenvector of covariance matrix which belongs to largest eigenvalue (which points in direction of largest variance)

↳ eigenvector belonging to 2<sup>nd</sup> largest eigenvalue points in direction of 2<sup>nd</sup> largest variance

General case: find M-dimensional principal subspace of D-dimensional dataset

$b_j, j = M+1, \dots, D$  (optimise these)

$$Sb_j = s_j b_j, \quad j = M+1, \dots, D$$

$$J = \sum_{j=M+1}^D s_j$$

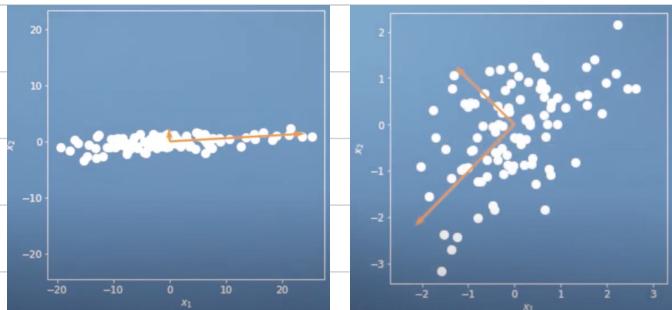
# PCA

Subtracting the mean can avoid numerical difficulties

↳ e.g. values centered around  $10^8$ : covariance matrix requires multiplication of large num., numerical instabilities

↳ then recommended to  $\div$  each dimension by corresponding  $\sigma$ : makes data unit free, variance in each dimension = 1

↳ does leave correlations intact



(unit variance)

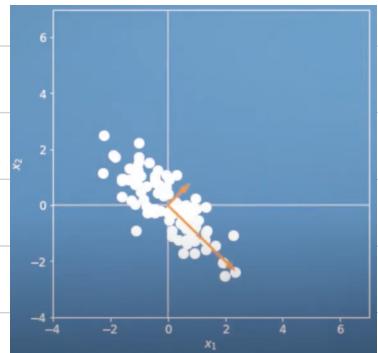
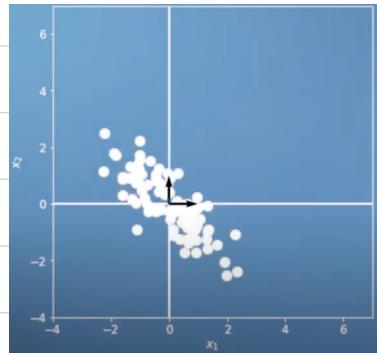
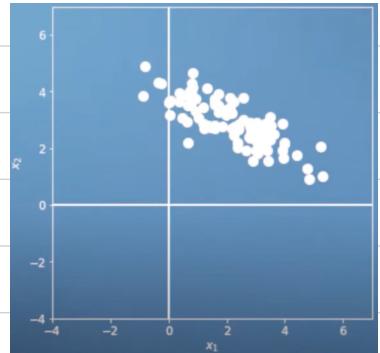
principal Subspace of normalised dataset

Strong correlation between 2 dimensions

\$x\_1\$ in cm, \$x\_2\$ in m

(So naturally varies more)

example:



eigenvectors are scaled by magnitude of corresponding eigenvalue  
longer vector Spans principal Subspace ( $U$ )  
Can then project any point  $x_*$  onto principal Subspace

(want to project onto 1D Subspace)

(center and  $\div \sigma$ , data is unit free with  $\sigma^2 = 1$  along each axis - black arrows)

(compute data covariance matrix and its eigenvalues and corresponding eigenvectors)

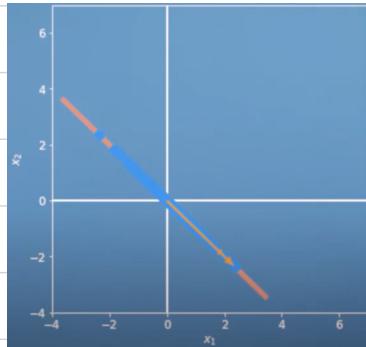
$$x_*^{(d)} \leftarrow \frac{x_*^{(d)} - \mu^{(d)}}{\sigma^{(d)}} \quad (\text{for every dimension in } x_*)$$

$$\tilde{x}_* = \pi_U(x_*) = B B^T x_*^{(d)}$$

↑  
projection of  $x_*$  onto principal Subspace  $U$

coordinates of projection w.r.t. basis of principal Subspace

matrix of eigenvectors belonging to largest eigenvalues as columns



projection of  $x^*$

Spanned by eigenvectors that belong to largest eigenvalues

## PCA in high dimensions

D-dimensions: covariance matrix is  $D \times D$  matrix

$\hookrightarrow$  high D: computing eigenvectors + eigenvalues is expensive (Scales cubically in no. of dimensions)

May have fewer data points than dimensions

$$x_1, \dots, x_N \in \mathbb{R}^D \quad (\text{assume centered data})$$

✓ i.e. Some dependencies

$$S = \frac{1}{N} X^T X \Rightarrow \text{rank}(S) = N, D-N+1 \text{ eigenvalues} = 0 \quad (\text{matrix not full rank, rows and cols. are linearly dependent})$$

$$X = \begin{bmatrix} x_1^T \\ \vdots \\ x_N^T \end{bmatrix} \in \mathbb{R}^{N \times D}$$

Turn  $D \times D$  covariance matrix  $S$  to full rank  $N \times N$  covariance matrix without eigenvalues = 0

$$S b_i = \lambda_i b_i \quad \text{basis vector of orthogonal complement of principal Subspace}$$

$$\underbrace{\frac{1}{N} X^T X}_{S} b_i = \lambda_i b_i$$

recover eigenvectors  
of  $S$

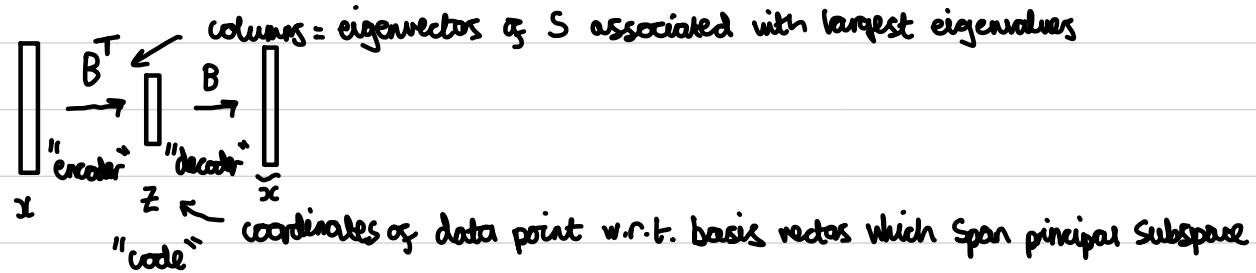
eigenvector of  $S$  belonging to eigenvalue  $\lambda_i$

$$\underbrace{\frac{1}{N} X X^T}_{\in \mathbb{R}^{N \times N}} \underbrace{X^T}_{C_i} b_i = \lambda_i \underbrace{X^T}_{C_i} b_i \quad (\text{multiply } X)$$

eigenvectors of

matrix (has same non-zero eigenvalues as  $S$ , but now  $N \times N$ , so can compute eigenvectors + eigenvalues much quicker  
than for original  $S$ )

## Perspectives of PCA



Find PCA parameters s.t. reconstruction error between  $x$  and  $\hat{x}$  is minimized

(can think of PCA as linear autoencoder (i.e. encode  $x$  and try to decode to something similar to  $x$ )

If encoder + decoder are linear mappings: get PCA solution by minimising squared autoencoding loss

↳ replace with non-linear mapping: non-linear autoencoder

↳ e.g. deep autoencoder where linear functions of encoder + decoder replaced with DNN's

Information theory perspective: code can be seen as smaller, compressed version of  $x$

↳ when reconstructing original data using code, don't get exact data back (slightly distorted / noisy version)

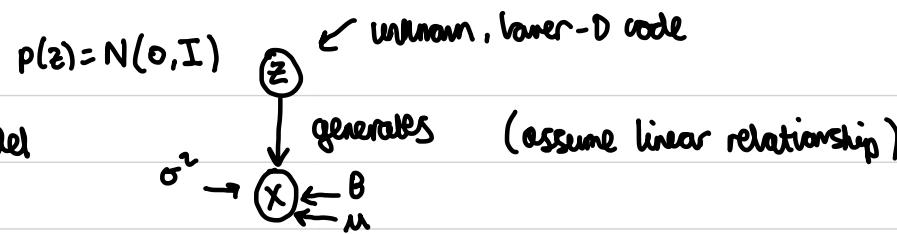
↳ lossy compression

↳ maximise correlation between original data and lower-D code

↳ related to mutual information (maximise)

Minimising variance = maximising variance of data when projected onto principal subspace

Interpret variance as info contained within data: PCA retains as much info as possible



(can look at PCA as latent variable model)

$$x = \beta z + \mu + \epsilon \quad \text{Likelihood: } p(x|z) = N(x|\beta z + \mu, \sigma^2 I)$$

$$\epsilon \sim N(0, \sigma^2 I) \quad \text{Marginal likelihood: } \int p(x|z) p(z) dz \\ = N(x|\mu, \beta\beta^T + \sigma^2 I)$$

Use maximum likelihood estimation to find parameters:

- $\mu$ : mean of data
- $\beta$ : matrix of eigenvectors corresponding to largest eigenvalues

To find  $z$  of data point: apply Bayes theorem to invert L.R. between  $z$  and  $x$

$$\hookrightarrow p(z|x) = \frac{p(x|z)p(z)}{p(x)}$$

5 different perspectives leading to different objectives:

- Minimising Squared reconstruction error
- Minimising autoencoder loss
- Maximising mutual information
- Maximising variance of projected data
- Maximising likelihood in latent variable model

(all give same solution to PCA problem)