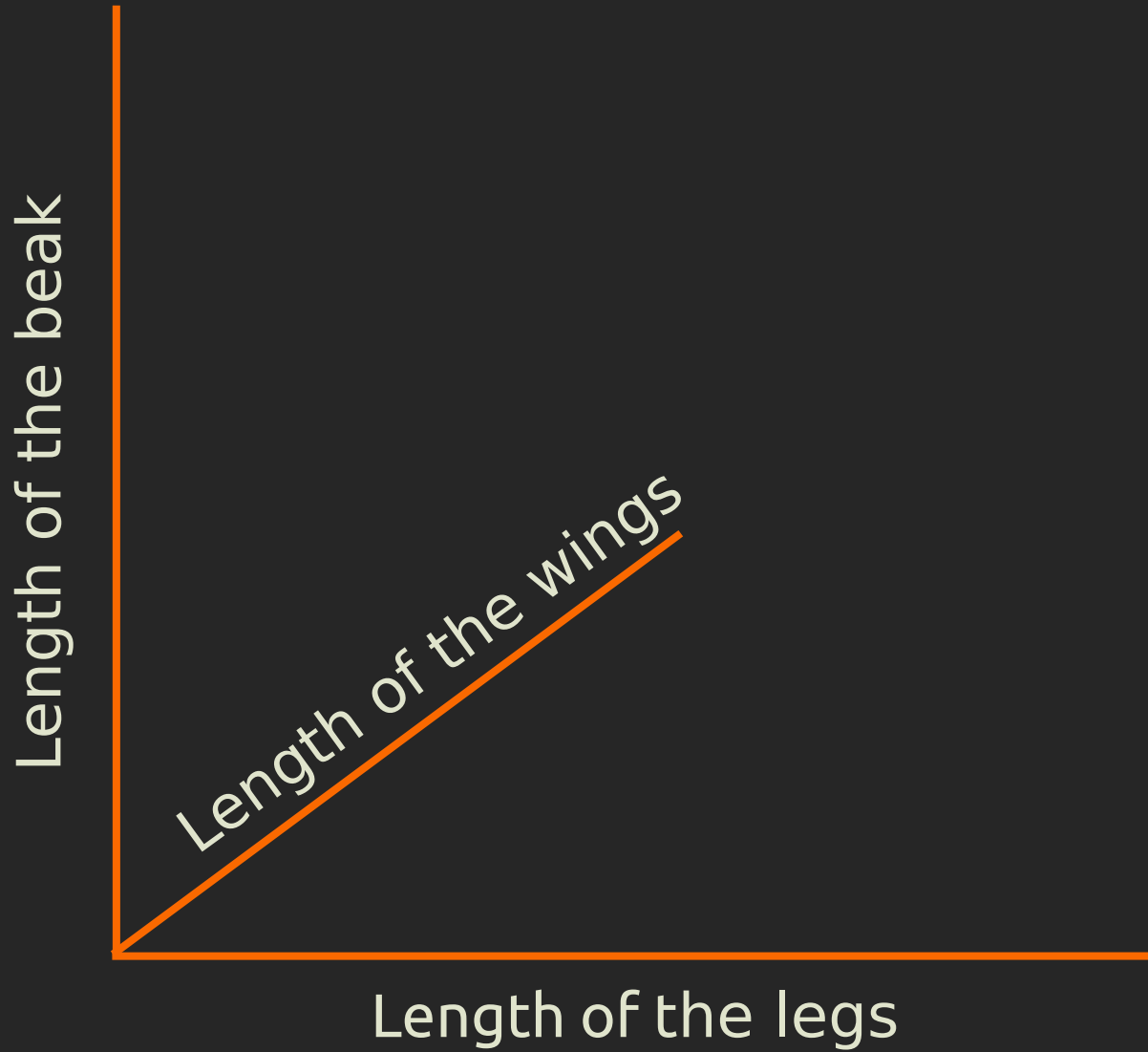


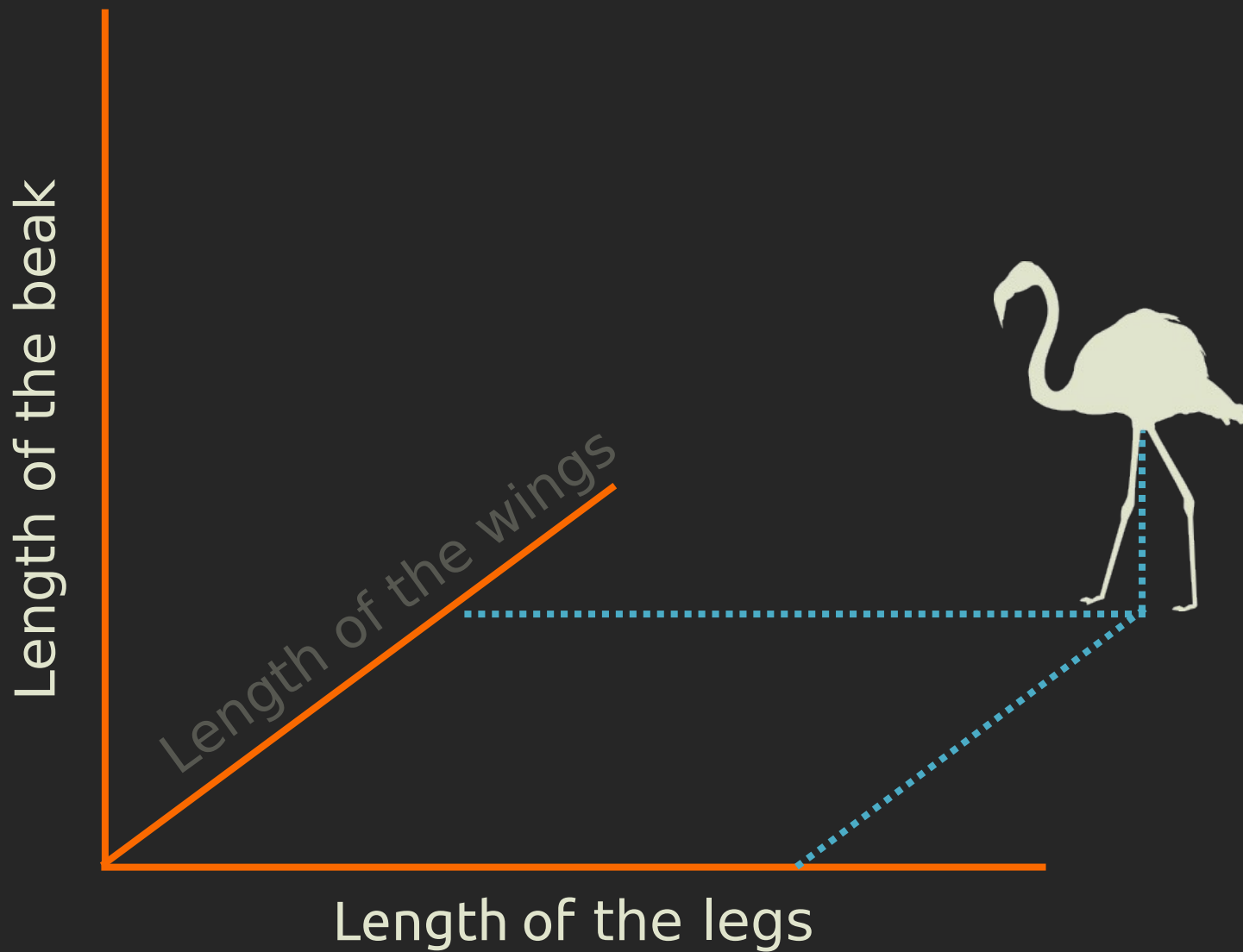
Ordinations

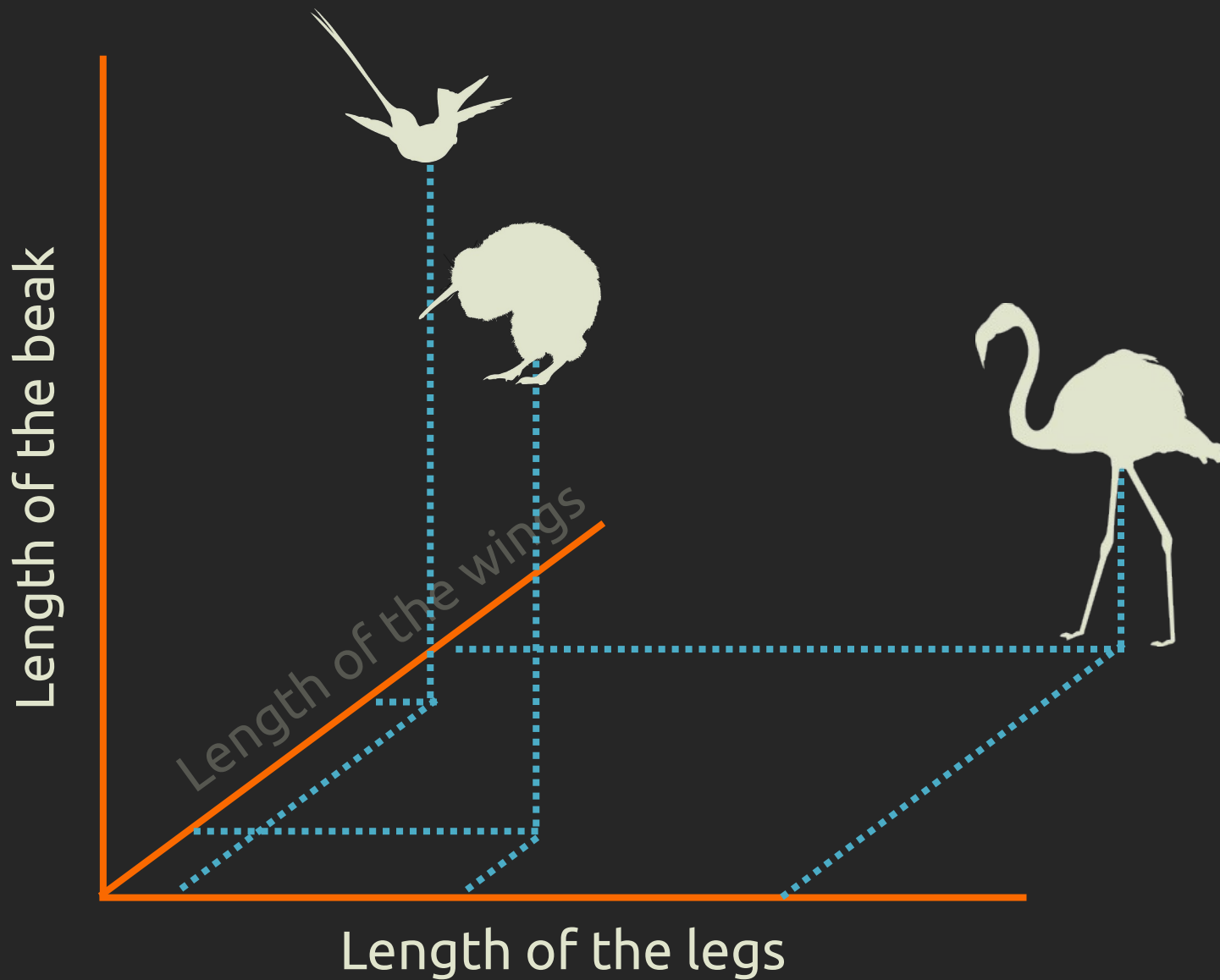
Length of the beak

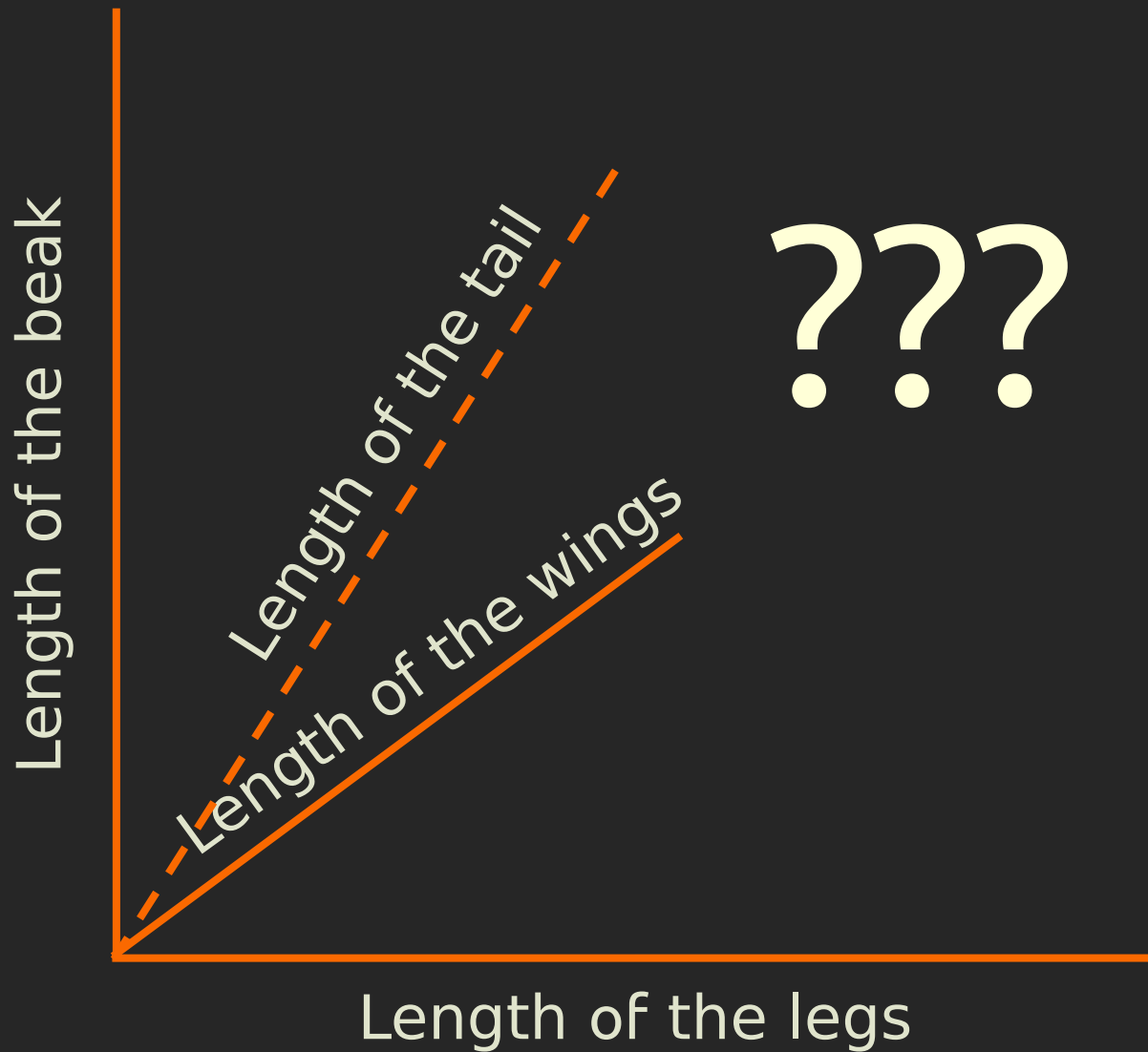


Length of the legs





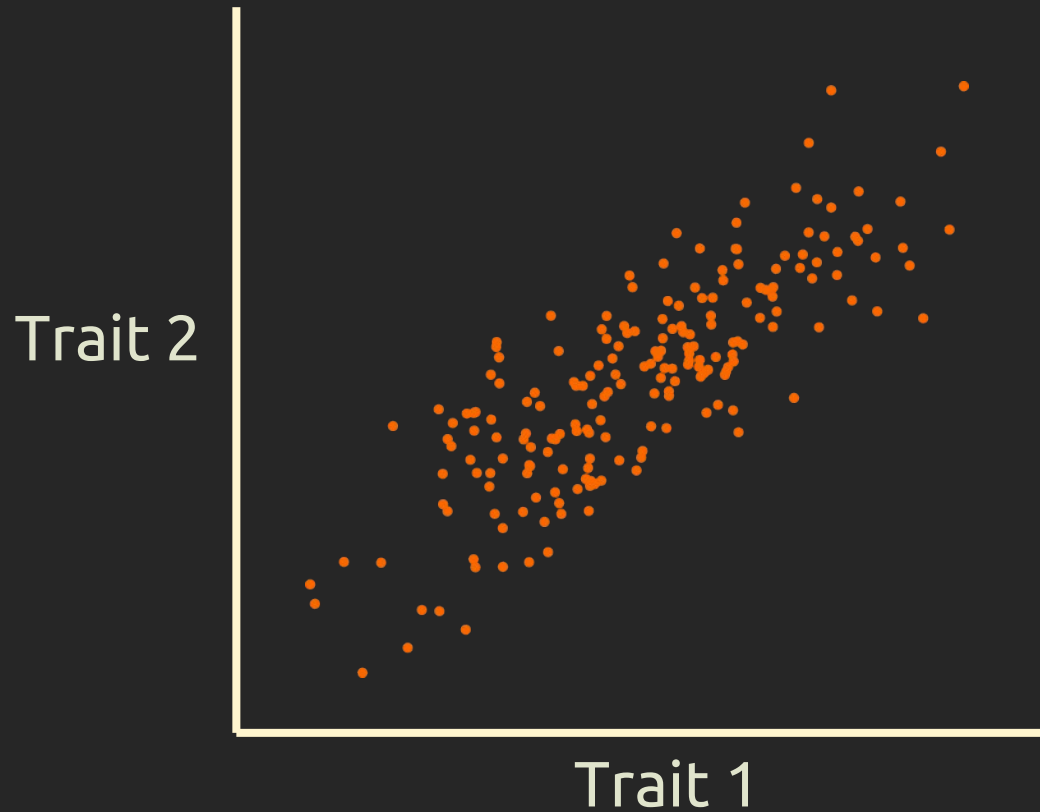




PCA

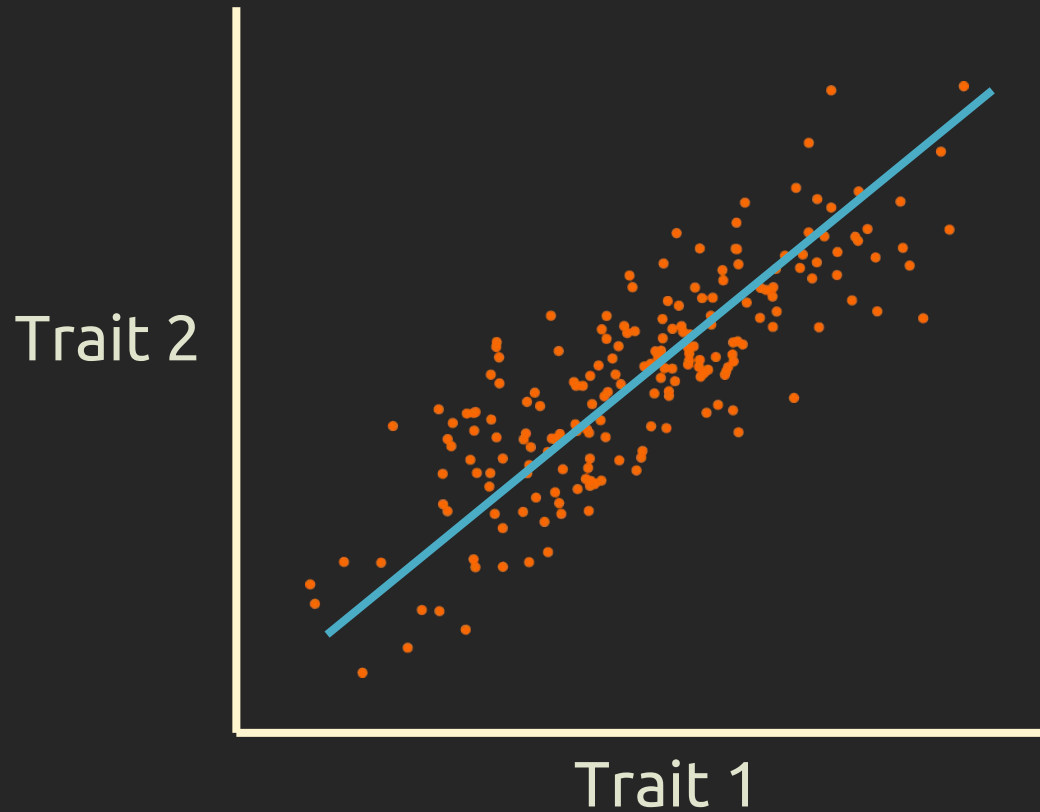
TL;DR:

PCA = rotating and stretching your data to
maximise variance covariance



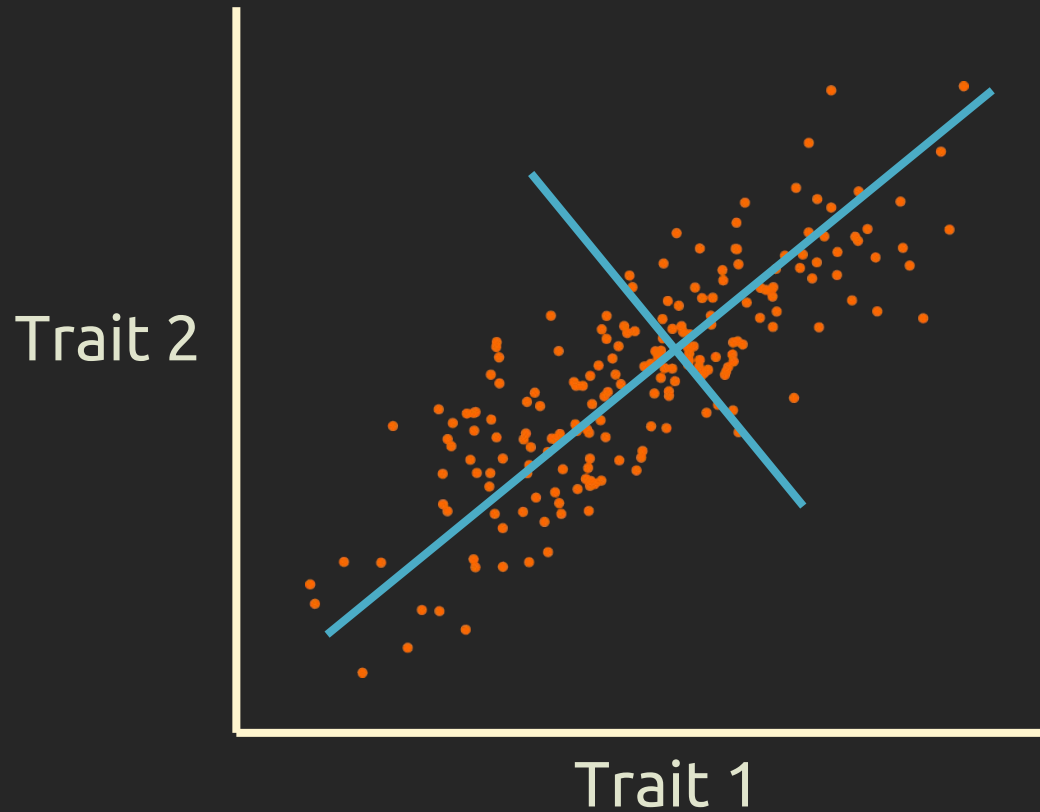
TL;DR:

PCA = rotating and stretching your data to
maximise variance covariance



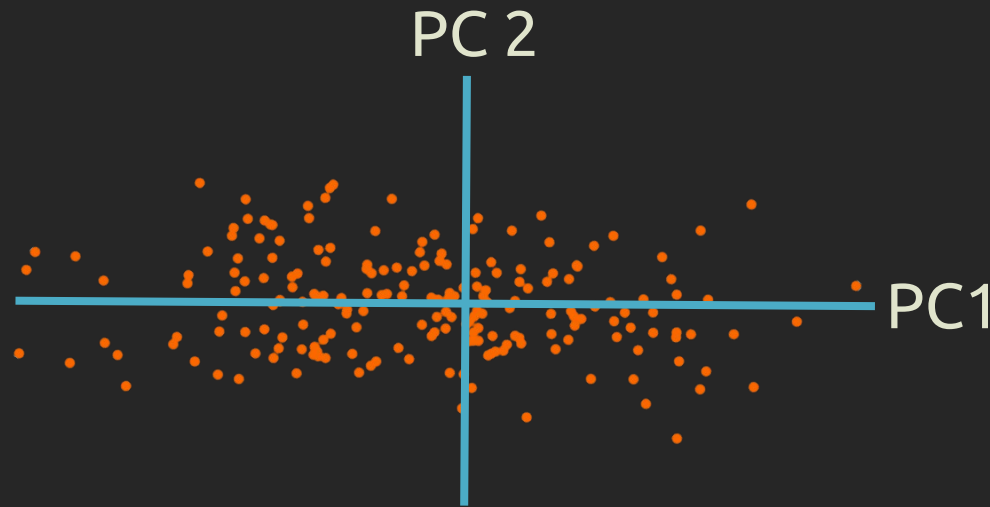
TL;DR:

PCA = rotating and stretching your data to maximise variance covariance



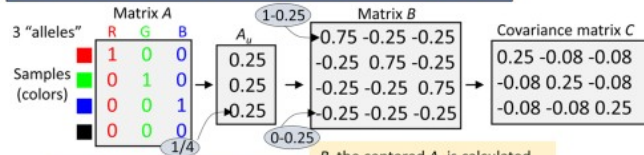
TL;DR:

PCA = rotating and stretching your data to
maximise variance covariance



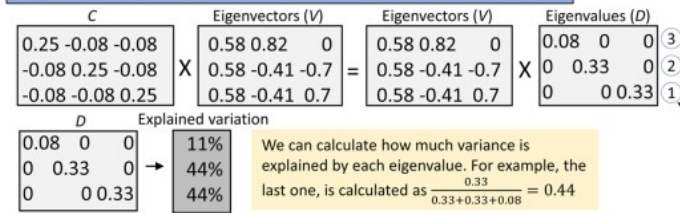
The black box

A. Step 1. Provided dataset A with four color samples and 3 "alleles", we will first calculate the covariance matrix C for the centered A matrix



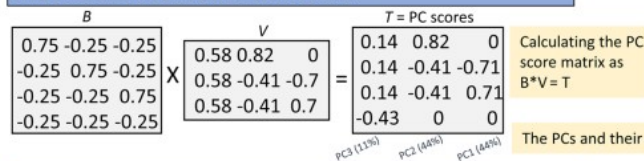
Calculate the covariance matrix for A as $C = \frac{1}{N-1} B^T B$ where $N=4$.

Step 2. Find the eigenvectors V and eigenvalues D for the covariance matrix C

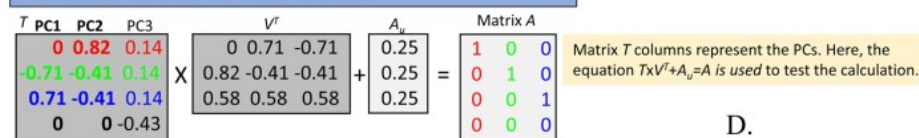


Find V and D that satisfy $Cx = Vx = Vx$. The eigenvalue columns are sorted by size, and the eigenvectors are sorted accordingly.

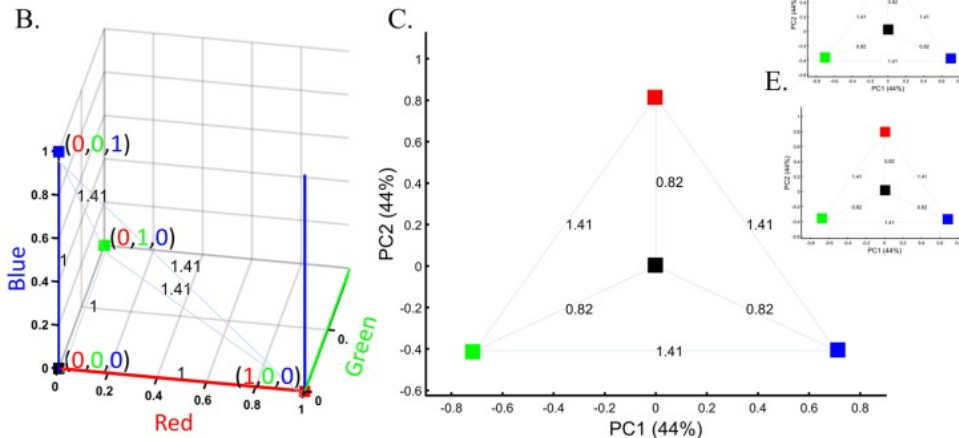
Step 3. Calculate the principal component score, T



Step 4. Validate that matrix A can be re-calculated

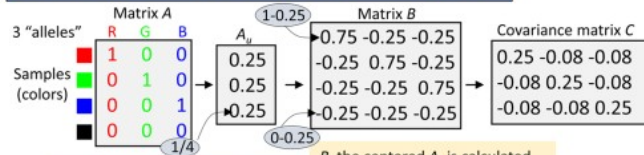


Step 5. Plot the top two PC scores as in C) below



Elhaik, E. Principal Component Analyses (PCA)-based findings in population genetic studies are highly biased and must be reevaluated. Sci Rep 12, 14683 (2022). <https://doi.org/10.1038/s41598-022-14395-4>

A. Step 1. Provided dataset A with four color samples and 3 "alleles", we will first calculate the covariance matrix C for the centered A matrix



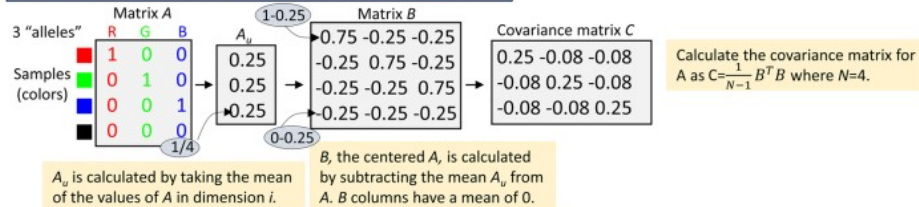
A_{ij} is calculated by taking the mean of the values of A in dimension i.

B, the centered A, is calculated by subtracting the mean A_{ij} from A. B columns have a mean of 0.

Step 1: **centre** the matrix

Step 2: measure the variance covariance of the centred matrix

A. Step 1. Provided dataset A with four color samples and 3 "alleles", we will first calculate the covariance matrix C for the centered A matrix

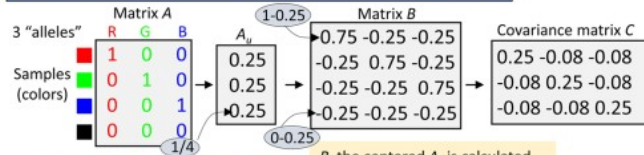


Keep this in a corner of your head, we'll get back to this!

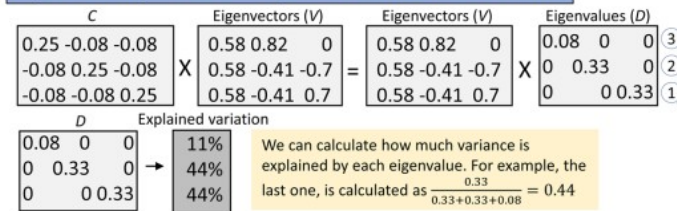
Step 1: **centre** the matrix

Step 2: measure the variance covariance of the centred matrix

A. Step 1. Provided dataset A with four color samples and 3 "alleles", we will first calculate the covariance matrix C for the centered A matrix



Step 2. Find the eigenvectors V and eigenvalues D for the covariance matrix C



Step 3: do a eigen decomposition.
Basically satisfy the equation:

VCV matrix * eigenvector = eigenvector * eigenvalue.

A. Step 1. Provided dataset A with four color samples and 3 "alleles", we will first calculate the covariance matrix C for the centered A matrix

3 "alleles" (R, G, B) and 4 "Samples (colors)" are used to calculate the covariance matrix C.

Matrix A (Allele counts):

	R	G	B
Sample 1 (Red)	1	0	0
Sample 2 (Green)	0	1	0
Sample 3 (Blue)	0	0	1
Sample 4 (Black)	0	0	0

Matrix A is converted to a 4x4 matrix A_{ij} by subtracting the mean (1/4) from each row:

	A_{11}	A_{12}	A_{13}	A_{14}
Sample 1	0.25	-0.25	-0.25	-0.25
Sample 2	-0.25	0.25	-0.25	-0.25
Sample 3	-0.25	-0.25	0.25	-0.25
Sample 4	-0.25	-0.25	-0.25	0.25

Matrix B (Standardized allele values):

	B_{11}	B_{12}	B_{13}	B_{14}
Sample 1	0.75	-0.25	-0.25	-0.25
Sample 2	-0.25	0.75	-0.25	-0.25
Sample 3	-0.25	-0.25	0.75	-0.25
Sample 4	-0.25	-0.25	-0.25	0.75

Covariance matrix C (Calculated from Matrix B):

	C_{11}	C_{12}	C_{13}
Allele R	0.25	-0.08	-0.08
Allele G	-0.08	0.25	-0.08
Allele B	-0.08	-0.08	0.25

The covariance matrix C is calculated as $C = \frac{1}{N-1} B^T B$, where $N=4$.

A_{ij} is calculated by taking the mean of the values of A in dimension i.

B, the centered A, is calculated by subtracting the mean A_{ij} from A. B columns have a mean of 0.

Step 2. Find the eigenvectors V and eigenvalues D for the covariance matrix C

C			Eigenvectors (V)			Eigenvectors (V)			Eigenvalues (D)		
0.25	-0.08	-0.08	0.58	0.82	0	0.58	0.82	0	0.08	0	0
-0.08	0.25	-0.08	0.58	-0.41	-0.7	0.58	-0.41	-0.7	0	0.33	0
-0.08	-0.08	0.25	0.58	-0.41	0.7	0.58	-0.41	0.7	0	0	0.33

Find V and D that satisfy $CxV = VxD$. The eigenvalue columns are sorted by size, and the eigenvectors are sorted accordingly.

D			Explained variation		
0.08	0	0	11%		
0	0.33	0	44%		
0	0	0.33	44%		

We can calculate how much variance is explained by each eigenvalue. For example, the last one, is calculated as $\frac{0.33}{0.33+0.33+0.08} = 0.44$

Step 3: do a eigen decomposition.
Basically satisfy the equation:

VCV matrix * **eigenvector** = **eigenvector** * **eigenvalue**.

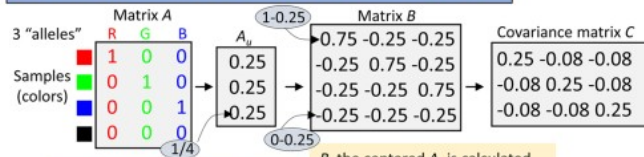
The variance-covariance between traits

The vector (direction) of change

The vector (direction) of change

The strength (length) of change

A. Step 1. Provided dataset A with four color samples and 3 "alleles", we will first calculate the covariance matrix C for the centered A matrix

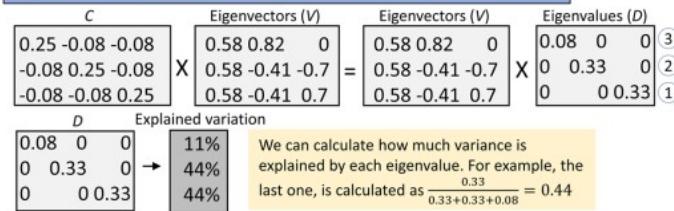


Calculate the covariance matrix for A as $C = \frac{1}{N-1} B^T B$ where $N=4$.

A_i is calculated by taking the mean of the values of A in dimension i .

B, the centered A, is calculated by subtracting the mean A_i from A. B columns have a mean of 0.

Step 2. Find the eigenvectors V and eigenvalues D for the covariance matrix C



Find V and D that satisfy $CxV = VxD$. The eigenvalue columns are sorted by size, and the eigenvectors are sorted accordingly.

We can calculate how much variance is explained by each eigenvalue. For example, the last one, is calculated as $\frac{0.33}{0.33+0.33+0.08} = 0.44$

Step 3: do a eigen decomposition. Basically satisfy the equation:

VCV matrix * eigenvector = eigenvector * eigenvalue.

This method is the core of the PCA. I think it's OK to treat it as a black box since it varies between algorithms. R default's is LAPACK but EISPACK or other algorithms can also be used. Also, these algorithms are explicitly approximations: "All you can hope for is a solution to a problem suitably close to x." (base::eigen). This can explain differences between ordinations of the same data.

A. Step 1. Provided dataset A with four color samples and 3 "alleles", we will first calculate the covariance matrix C for the centered A matrix

Matrix A				Matrix B				Covariance matrix C			
3 "alleles"	R	G	B								
1	1	0	0	0.25	0.75	-0.25	-0.25	0.25	-0.08	-0.08	
2	0	1	0	0.25	-0.25	0.75	-0.25	-0.08	0.25	-0.08	
3	0	0	1	0.25	-0.25	-0.25	0.75	-0.08	-0.08	0.25	
4	0	0	0	0.25	-0.25	-0.25	-0.25	-0.08	-0.08	0.25	

A_{ij} is calculated by taking the mean of the values of A in dimension i.

B, the centered A, is calculated by subtracting the mean A_{ij} from A. B columns have a mean of 0.

Calculate the covariance matrix for A as $C = \frac{1}{N-1} B^T B$ where $N=4$.

Step 2. Find the eigenvectors V and eigenvalues D for the covariance matrix C

C	Eigenvectors (V)	Eigenvectors (V)	Eigenvalues (D)
0.25 -0.08 -0.08 -0.08 0.25 -0.08 -0.08 -0.08 0.25	0.58 0.82 0 0.58 -0.41 -0.7 0.58 -0.41 0.7	0.58 0.82 0 0.58 -0.41 -0.7 0.58 -0.41 0.7	0.08 0 0 0 0.33 0 0 0 0.33
D	Explained variation		
0.08 0 0 0 0.33 0 0 0 0.33	11% 44% 44%		

We can calculate how much variance is explained by each eigenvalue. For example, the last one, is calculated as $\frac{0.33}{0.33+0.33+0.08} = 0.44$

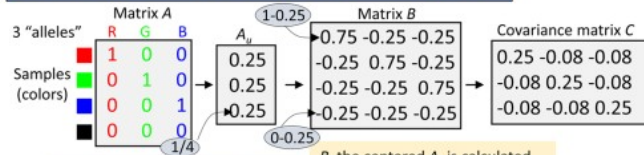
Find V and D that satisfy $CxV = VxD$. The eigenvalue columns are sorted by size, and the eigenvectors are sorted accordingly.

Step 3. Calculate the principal component score, T

B	V	T = PC scores	
0.75 -0.25 -0.25 -0.25 0.75 -0.25 -0.25 -0.25 0.75 -0.25 -0.25 -0.25	0.58 0.82 0 0.58 -0.41 -0.7 0.58 -0.41 0.7	0.14 0.82 0 0.14 -0.41 -0.71 0.14 -0.41 0.71 -0.43 0 0	Calculating the PC score matrix as $B \cdot V = T$
			The PCs and their explained variation (%)
			PC3 (11%) PC2 (44%) PC1 (44%)

Step 4: multiply the centred matrix by the eigenvector

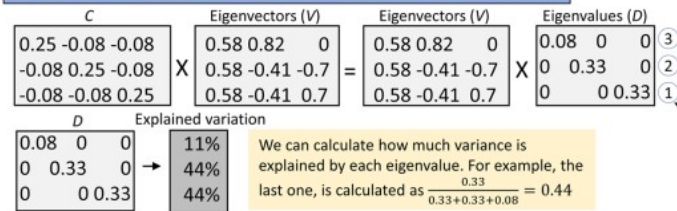
A. Step 1. Provided dataset A with four color samples and 3 "alleles", we will first calculate the covariance matrix C for the centered A matrix



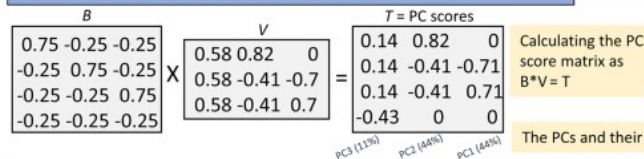
A_i is calculated by taking the mean of the values of A in dimension i .

B, the centered A, is calculated by subtracting the mean A_i from A. B columns have a mean of 0.

Step 2. Find the eigenvectors V and eigenvalues D for the covariance matrix C

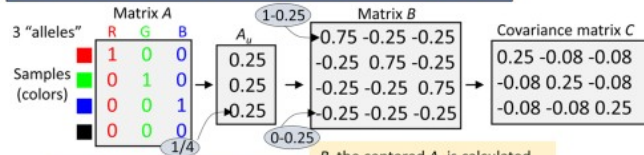


Step 3. Calculate the principal component score, T



Step 4: multiply the centred matrix by the eigenvector
Step 5: that's it.

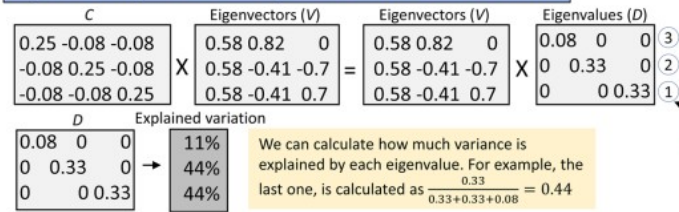
A. Step 1. Provided dataset A with four color samples and 3 "alleles", we will first calculate the covariance matrix C for the centered A matrix



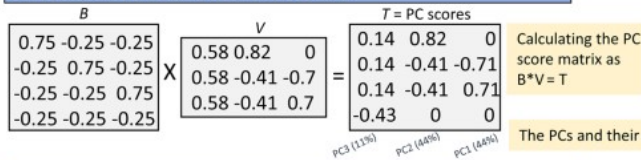
A_u is calculated by taking the mean of the values of A in dimension i.

B, the centered A, is calculated by subtracting the mean A_u from A. B columns have a mean of 0.

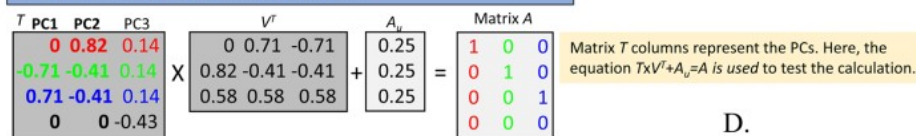
Step 2. Find the eigenvectors V and eigenvalues D for the covariance matrix C



Step 3. Calculate the principal component score, T

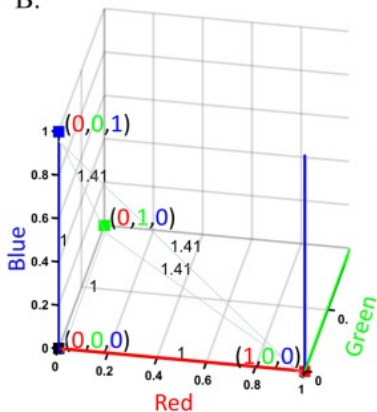


Step 4. Validate that matrix A can be re-calculated

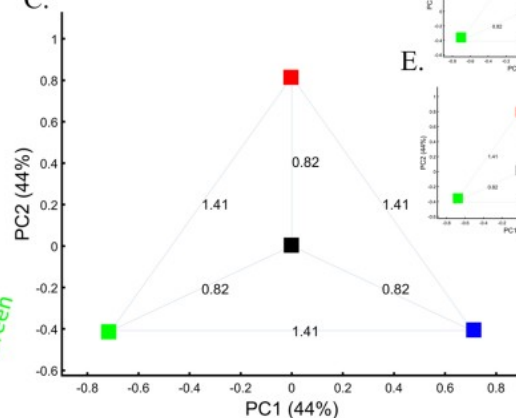


Step 5. Plot the top two PC scores as in C) below

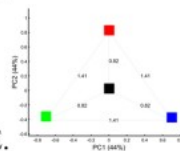
B.



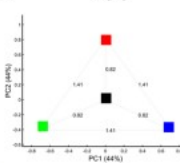
C.



D.



E.



PCA: what it's good at?

- *Ordinating* your data (i.e. ranking all your variables and making them independent and orthogonal).

PCA: what it's good at?

- *Ordinating* your data (i.e. ranking all your variables and making them independent and orthogonal).
- Reducing dimensionality (to some extent – e.g. going from 200D to 10D).

PCA: what it's good at?

- *Ordinating* your data (i.e. ranking all your variables and making them independent and orthogonal).
- Reducing dimensionality (to some extent – e.g. going from 200D to 10D).
- Creating a “true” mathematical space (that contains all the possible trait combinations).

PCA: what it's **bad** at (in my opinion)?

- Being interpreted by humans in 2D.

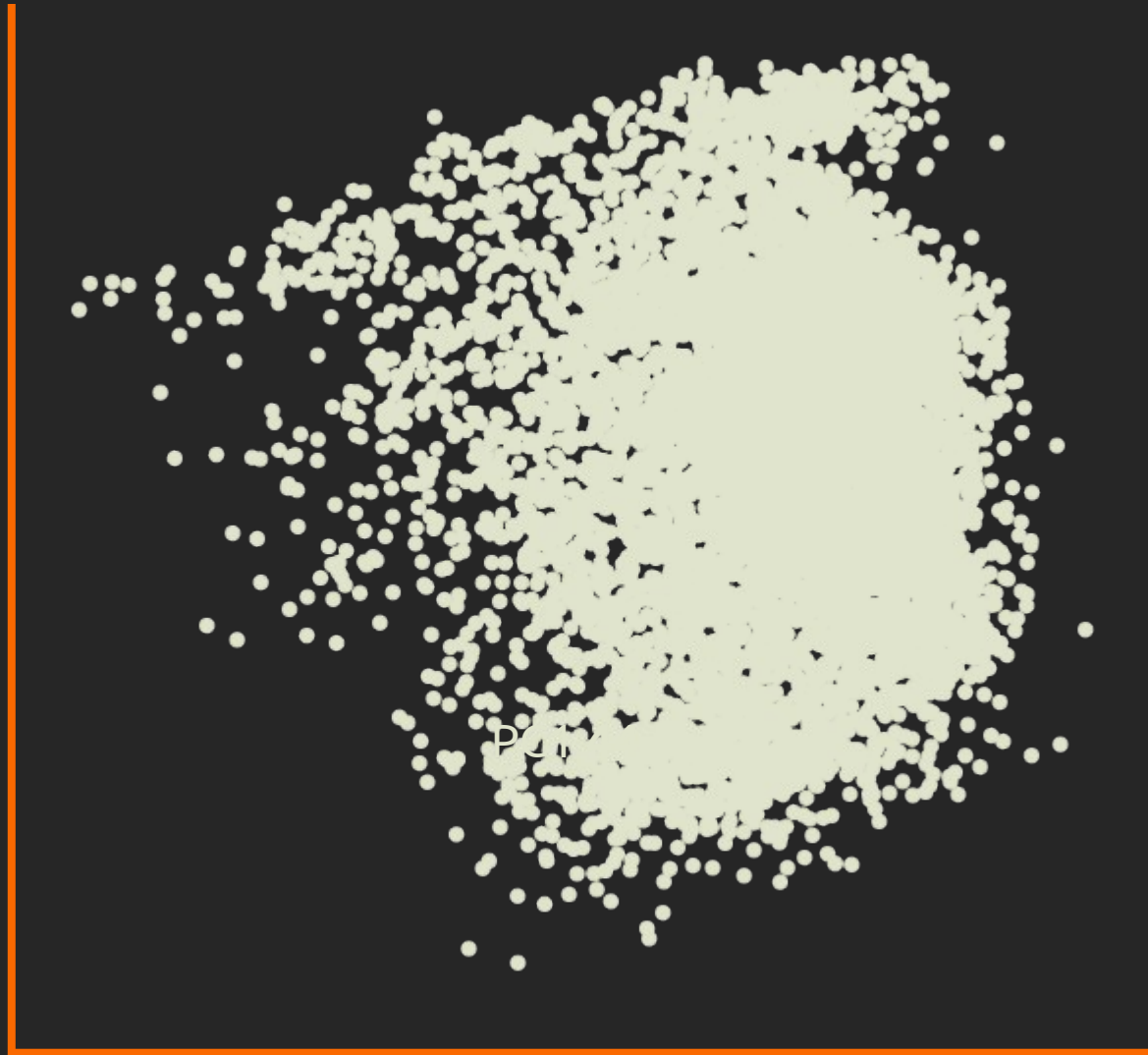
PCA: what it's **bad** at (in my opinion)?

- Being interpreted by humans in 2D.
- Spreading data (the centre of your space is the average data value, not anything biological).

PCA: what it's **bad** at (in my opinion)?

- Being interpreted by humans in 2D.
- Spreading data (the centre of your space is the average data value, not anything biological).
- Creating dimensions that are easy to interpret (e.g. PC1 = correlation between n-variables decided by the algorithm – these can sometimes map to biological things, sometimes not!).

PC2
5.93%



PC1 89.19%

