

Principle Component Analysis (PCA)

Dr Jamie A Ward

This week

Vector spaces

- ▶ Linear dependence
- ▶ Mean and variance
- ▶ Covariance

Principle Component Analysis (PCA)

- ▶ Maximise the covariance
- ▶ Linearly independent features

Linearly dependent vectors

How can you tell if two vectors are linearly dependent or independent?

$$x = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} \quad y = \begin{bmatrix} y_1 \\ y_2 \\ y_3 \end{bmatrix}$$

If one can be written as a **linear combination** of the other, they are linearly **dependent**. Otherwise, they are linearly independent.

x and y are linearly dependent *iff* there exists scalars a, b such that:

$$a \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = b \begin{bmatrix} y_1 \\ y_2 \\ y_3 \end{bmatrix} \quad \text{or equivalently,} \quad \begin{bmatrix} x_1 & y_1 \\ x_2 & y_2 \\ x_3 & y_3 \end{bmatrix} \begin{bmatrix} a \\ b \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix}$$

Otherwise, they are linearly independent.

Orthogonal vectors

How can you tell if two vectors are orthogonal (at right angles) to each other?

$$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} \quad \mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \\ y_3 \end{bmatrix}$$

If the dot product of two (non-zero) vectors is 0, they are orthogonal:

$$\mathbf{x} \cdot \mathbf{y} = \mathbf{x}^T \mathbf{y} = [x_1 \ x_2 \ x_3] \begin{bmatrix} y_1 \\ y_2 \\ y_3 \end{bmatrix} = x_1 y_1 + x_2 y_2 + x_3 y_3 = 0$$

e.g.

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

Orthogonal vectors

How can you tell if two vectors are orthogonal (at right angles) to each other?

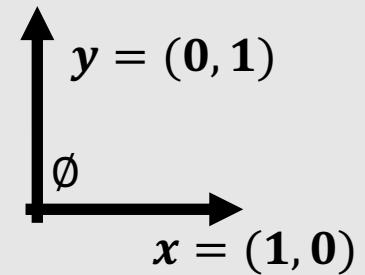
$$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} \quad \mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \\ y_3 \end{bmatrix}$$

If the dot product of two (non-zero) vectors is 0, they are orthogonal:

$$\mathbf{x} \cdot \mathbf{y} = \mathbf{x}^T \mathbf{y} = [x_1 \ x_2 \ x_3] \begin{bmatrix} y_1 \\ y_2 \\ y_3 \end{bmatrix} = x_1 y_1 + x_2 y_2 + x_3 y_3 = 0$$

This is because the dot product is proportional to the **cosine** of the angle θ between \mathbf{x} and \mathbf{y}

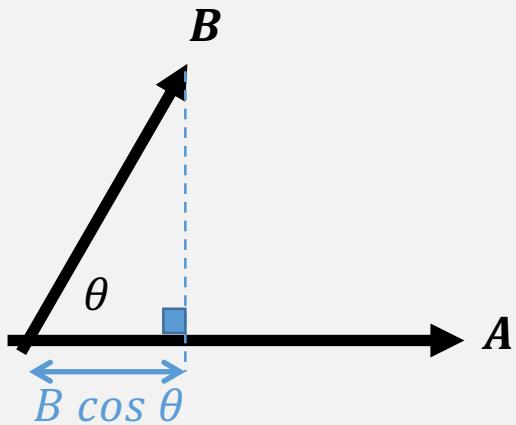
$$\mathbf{y} = \begin{bmatrix} 0 \\ 1 \end{bmatrix}, \mathbf{x} = \begin{bmatrix} 1 \\ 0 \end{bmatrix} \Rightarrow \mathbf{x} \cdot \mathbf{y} \propto \cos(\theta) = \cos(90^\circ) = 0$$



Aside: Geometry of the dot product

The dot product of two vectors is the magnitude of one times the projection of the second onto the first.

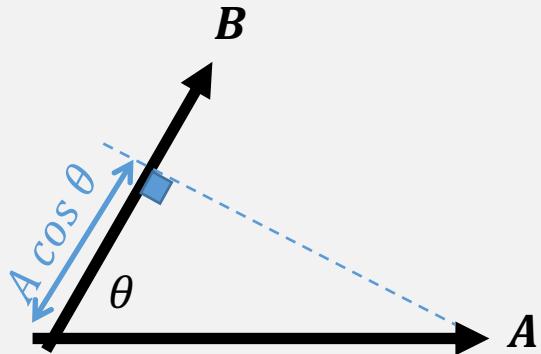
$$\mathbf{A} \cdot \mathbf{B} = AB \cos \theta$$



Aside: Geometry of the dot product

The dot product of two vectors is the magnitude of one times the projection of the second onto the first.

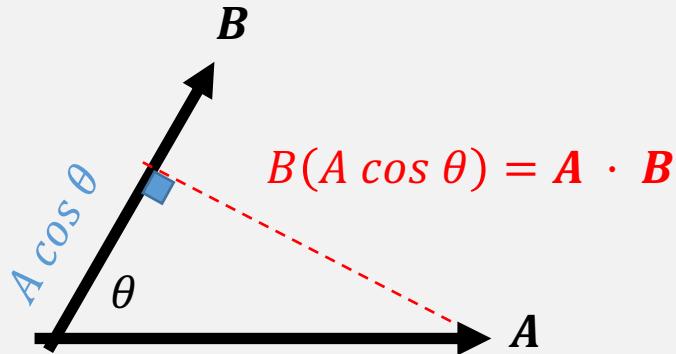
$$\mathbf{A} \cdot \mathbf{B} = AB \cos \theta$$



Aside: Geometry of the dot product

The dot product of two vectors is the magnitude of one times the projection of the second onto the first.

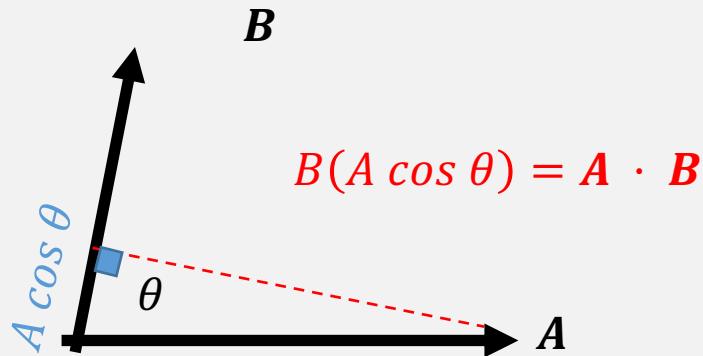
$$\mathbf{A} \cdot \mathbf{B} = AB \cos \theta$$



Aside: Geometry of the dot product

The dot product of two vectors is the magnitude of one times the projection of the second onto the first.

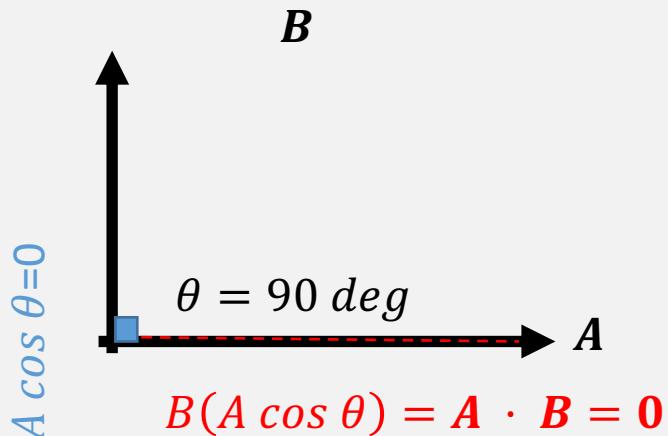
$$\mathbf{A} \cdot \mathbf{B} = AB \cos \theta$$



Aside: Geometry of the dot product

The dot product of two vectors is the magnitude of one times the projection of the second onto the first.

$$\mathbf{A} \cdot \mathbf{B} = AB \cos \theta$$



Linear dependence and orthogonality

Given vectors \mathbf{x} and \mathbf{y}

Two tests

1. If dot product $\mathbf{x} \cdot \mathbf{y} = 0$, then the vectors are orthogonal and thus **linearly independent**.
2. They are **linearly dependent** if there exists non-zero scalars, a and b , that satisfy:

$$a \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = b \begin{bmatrix} y_1 \\ y_2 \\ y_3 \end{bmatrix} \quad \text{or equivalently,} \quad \begin{bmatrix} x_1 & y_1 \\ x_2 & y_2 \\ x_3 & y_3 \end{bmatrix} \begin{bmatrix} a \\ b \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix}$$

Otherwise the vectors are **linearly independent**.

Linear dependence question

- Is x linearly dependent on y ?

$$x = \begin{bmatrix} 2 \\ 4 \\ 0 \\ 1 \end{bmatrix} \quad y = \begin{bmatrix} 4 \\ 8 \\ 0 \\ 2 \end{bmatrix}$$

Yes, since $ay = bx$, $a = 2, b = 1$

$$x = \begin{bmatrix} 0 \\ 4 \\ 0 \end{bmatrix} \quad y = \begin{bmatrix} 2 \\ 0 \\ 1 \end{bmatrix}$$

No, since $x^T y = 0$, i.e.

$$x^T y = [0 \ 4 \ 0] \begin{bmatrix} 2 \\ 0 \\ 1 \end{bmatrix} = 0 * 2 + 4 * 0 + 0 * 1 = 0$$

Or, equivalently, $\nexists \{a, b\}$ s.t. $ay = bx$

$$x = \begin{bmatrix} 1 \\ 2 \end{bmatrix} \quad y = \begin{bmatrix} 2 \\ 1 \end{bmatrix}$$

No, since $\nexists \{a, b\}$ s.t. $ay = bx$

\nexists means does not exist

Sample mean and variance

Given a random variable X that takes on a set of N values $x^{(i)}$, then the expected value of X is defined as:

$$E[x] = \sum_{i=1}^N x^{(i)} p_i$$

(where p_i is the *probability* of x_i appearing)

If we assume $p_i = \frac{1}{N}$, then:

$$E[x] = \text{mean}(x) = \bar{x} = \frac{1}{N} \sum_{i=1}^N x^{(i)}$$

The variance of X can be written:

$$\text{var}(x) = \sigma^2 = \frac{1}{N} \sum_{i=1}^N (x^{(i)} - \bar{x})^2$$

Sample mean and variance: example

You are given a dataset with the scalar values:

$$X = [0,1,2,3,4]$$

Estimate the mean as:

$$\bar{x} = \frac{1}{N} \sum_{i=1}^N x^{(i)} = \frac{0 + 1 + 2 + 3 + 4}{5} = 2$$

Using this, the variance of X can be estimated as:

$$\sigma^2 = \frac{1}{N} \sum_{i=1}^N (x^{(i)} - \bar{x})^2 = \frac{(-2)^2 + (-1)^2 + 0 + 1 + 2^2}{5} = \frac{10}{5} = 2$$

How does \bar{x} and σ^2 change if we grow the dataset as follows (does it increase, decrease, or stay the same)?

1. $X = [0,1,2,3,4,2,2,2,2]$ $\bar{x} = 2$ $\sigma^2 = 1.11$

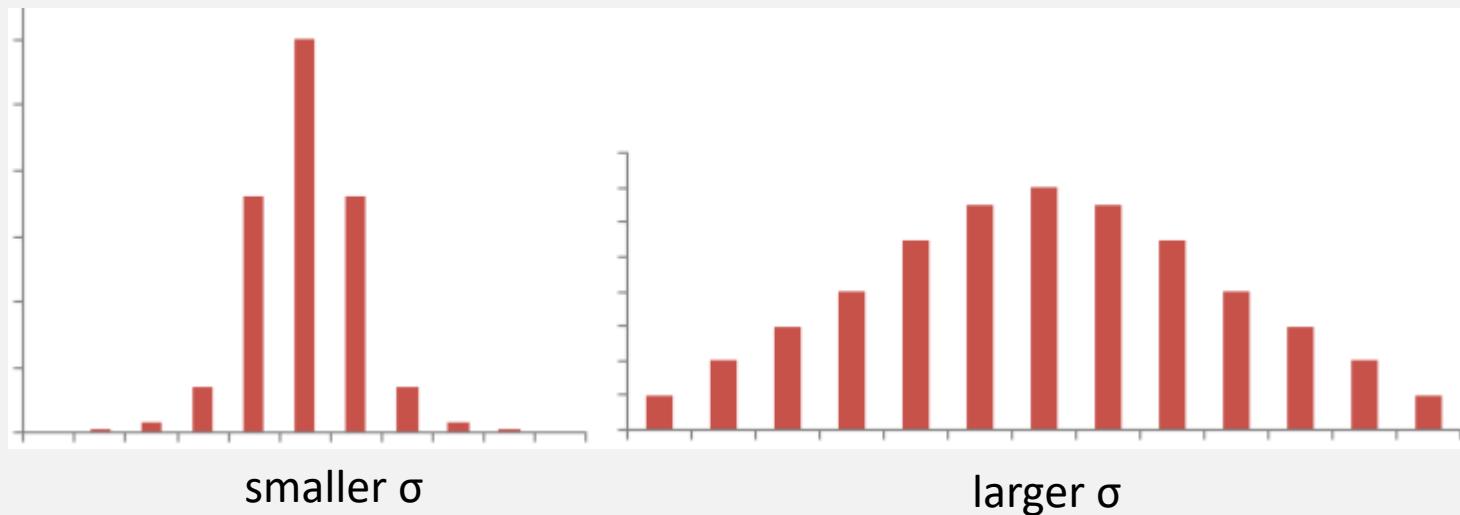
2. $X = [0,1,2,3,4, -1,5]$ $\bar{x} = 2$ $\sigma^2 = 4$

Variance: how far data is from the mean

$$X = [0,1,2,3,4] \quad \bar{x} = 2 \quad \sigma^2 = 2$$

$$X = [0,1,2,3,4,2,2,2,2] \quad \bar{x} = 2 \quad \sigma^2 = 1.11$$

$$X = [0,1,2,3,4, -1,5] \quad \bar{x} = 2 \quad \sigma^2 = 4$$



Mean and Covariance over >1 features

You are given a data matrix \mathbf{X} with F features and N samples:

$$\mathbf{X} = [\mathbf{x}^{(1)} \mathbf{x}^{(2)}, \dots \mathbf{x}^{(N)}], \quad \text{where } \mathbf{X} \in \mathbb{R}^{F \times N}$$

The sample mean is then:

$$\bar{\mathbf{x}} = \frac{1}{N} \sum_{i=1}^N \mathbf{x}^{(i)}, \quad \text{where } \mathbf{x}^{(i)} \in \mathbb{R}^{F \times 1}$$

The covariance is estimated as:

$$cov(\mathbf{X}) = \mathbf{S} = \frac{1}{N} \sum_{i=1}^N (\mathbf{x}^{(i)} - \bar{\mathbf{x}})(\mathbf{x}^{(i)} - \bar{\mathbf{x}})^T$$

This captures how the N variables vary together – how they change together with respect to the mean.

Mean, variance, and covariance

1 Dimension

$$\begin{aligned} \mathbf{X} &= [0,1,2,3,4] && \text{Data} \\ \bar{x} &= \frac{1}{N} \sum_{i=1}^N x^{(i)} = 2 && \text{Mean (scalar)} \\ \sigma^2 &= \frac{1}{N} \sum_{i=1}^N (x^{(i)} - \bar{x})^2 = 2 && \text{Variance (scalar)} \end{aligned}$$

> 1 dimension (x is an F-dimensional vector)

$$\begin{aligned} \bar{\mathbf{x}} &= \frac{1}{N} \sum_{i=1}^N \mathbf{x}^{(i)} && \text{Mean (Fx1)} \\ cov(\mathbf{X}) = \mathbf{S} &= \frac{1}{N} \sum_{i=1}^N (\mathbf{x}^{(i)} - \bar{\mathbf{x}})(\mathbf{x}^{(i)} - \bar{\mathbf{x}})^T && \text{Covariance (FxF)} \end{aligned}$$

Covariance

- How two feature vectors change together is based on the product of two scalars, e.g. s1 and s2:

S1	S2	S1*S2
-	-	+
+	+	+
-	+	-
+	-	-

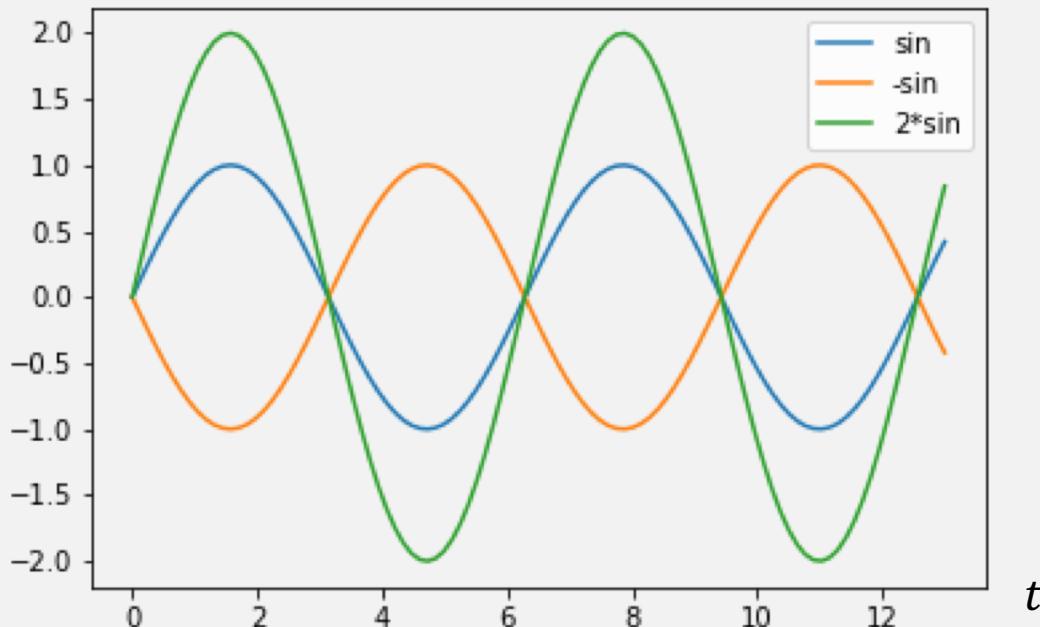
+ve sign: Covariance increase

-ve sign: Covariance decrease

Covariance example

Given data with N=100 samples,
F=3 features, and zero mean:

feature1 $x_1 = \sin(t)$
feature2 $x_2 = -\sin(t)$
feature3 $x_3 = 2 * \sin(t)$



1. What are the dimensions of feature matrix X ?
2. What are the dimensions of covariance matrix S ?
3. Is the covariance of feature1 to feature3 positive or negative?
4. What about the covariance of feature1 to feature2?

Covariance example

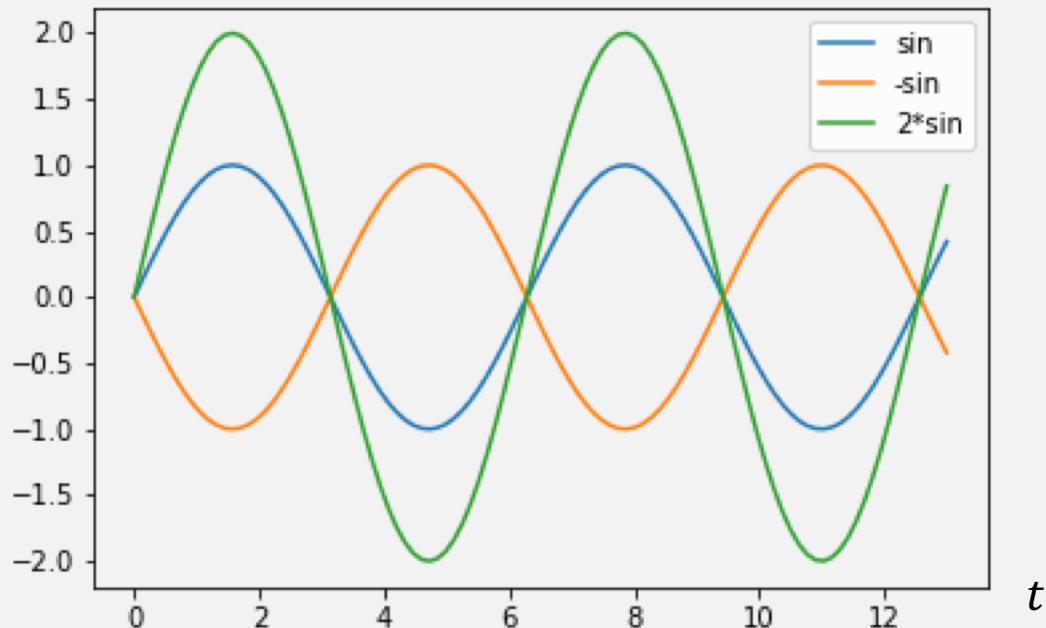
N=100 samples, F=3 features
 $X \in \mathbb{R}^{F \times N}, F = 3, N = 100$

Dimensions of S are 3×3
(Features x Features)

$$S = \frac{1}{N} \sum_{i=1}^N (x^{(i)} - \bar{x})(x^{(i)} - \bar{x})^T$$

$$S = \begin{bmatrix} 0.5 & -0.5 & 1 \\ -0.5 & 0.5 & -1 \\ 1 & -1 & 2 \end{bmatrix}$$

feature1 $x_1 = \sin(t)$
feature2 $x_2 = -\sin(t)$
feature3 $x_3 = 2 * \sin(t)$



1. What are the dimensions of feature matrix X ?
2. What are the dimensions of covariance matrix S ?
3. Is the covariance of feature1 to feature3 positive or negative?
4. What about the covariance of feature1 to feature2?

Covariance example

N=100 samples, F=3 features
 $X \in \mathbb{R}^{F \times N}, F = 3, N = 100$

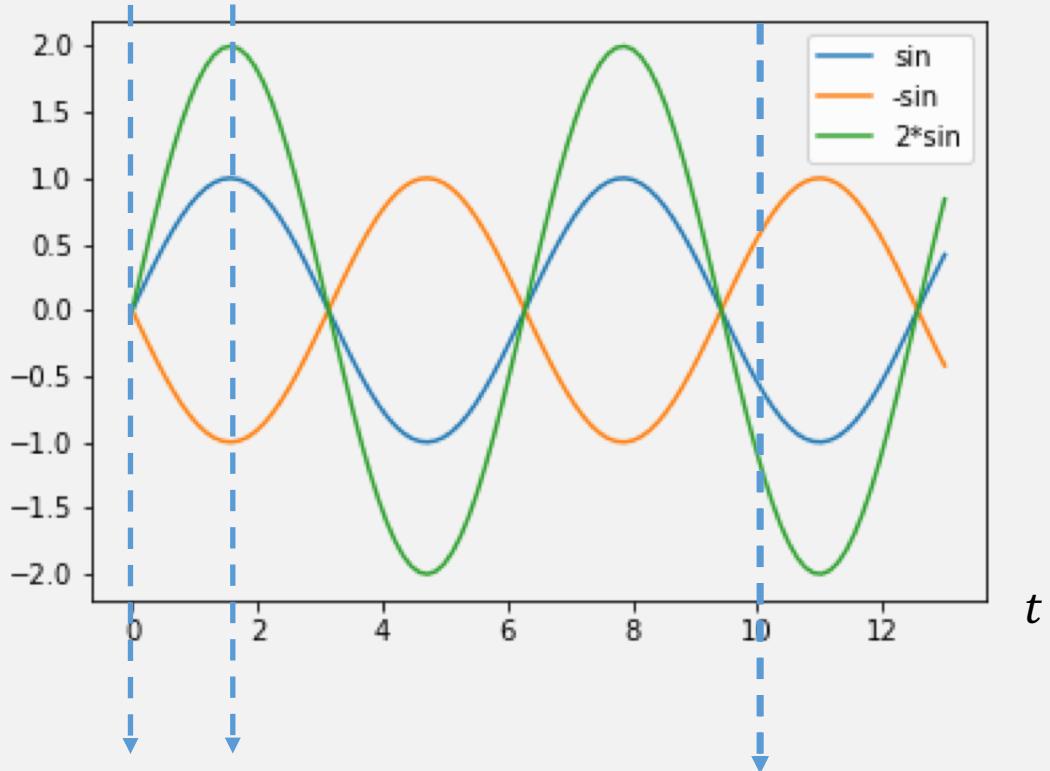
Dimensions of S are 3×3
(Features x Features)

$$S = \frac{1}{N} \sum_{i=1}^N (x^{(i)} - \bar{x})(x^{(i)} - \bar{x})^T$$

$$S = \begin{bmatrix} 0.5 & -0.5 & 1 \\ -0.5 & 0.5 & -1 \\ 1 & -1 & 2 \end{bmatrix}$$

$$\boldsymbol{x}^{(0)} = \begin{bmatrix} x_1^{(0)} & x_2^{(0)} & x_3^{(0)} \end{bmatrix}^T = \begin{bmatrix} 0 \\ 0 \\ 0 \end{bmatrix} \quad \boldsymbol{x}^{(12)} = \begin{bmatrix} 1 \\ -1 \\ 2 \end{bmatrix} \quad \boldsymbol{x}^{(76)} = \begin{bmatrix} -0.53 \\ 0.52 \\ -1.05 \end{bmatrix}$$

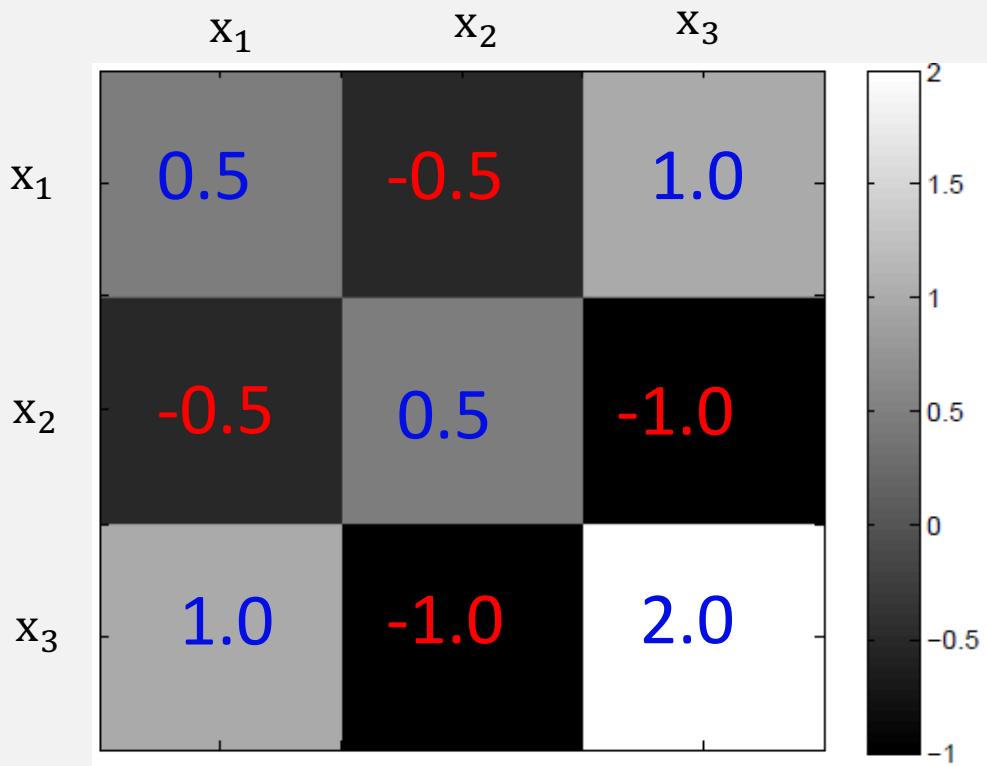
feature1 $x_1 = \sin(t)$
feature2 $x_2 = -\sin(t)$
feature3 $x_3 = 2 * \sin(t)$



Covariance example

$$S = \frac{1}{N} \sum_{i=1}^N (x^{(i)} - \bar{x})(x^{(i)} - \bar{x})^T$$

$$S = \begin{bmatrix} 0.5 & -0.5 & 1 \\ -0.5 & 0.5 & -1 \\ 1 & -1 & 2 \end{bmatrix}$$



$\text{Cov}(x_1, x_3) = \text{Cov}(x_3, x_1) = 1 \rightarrow +\text{ve}: x_1 \text{ and } x_3 \text{ move together}$

$\text{Cov}(x_1, x_2) = \text{Cov}(x_2, x_1) = -0.5 \rightarrow -\text{ve}: x_1 \text{ and } x_2 \text{ move apart}$

This week

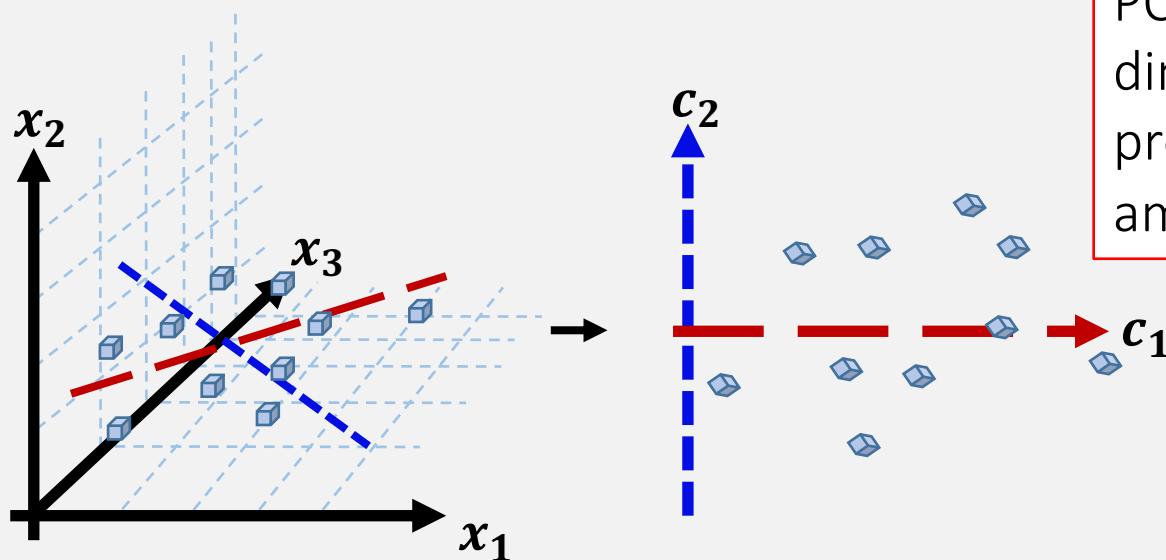
Vector spaces

- ▶ Linear dependence
- ▶ Mean and variance
- ▶ Covariance

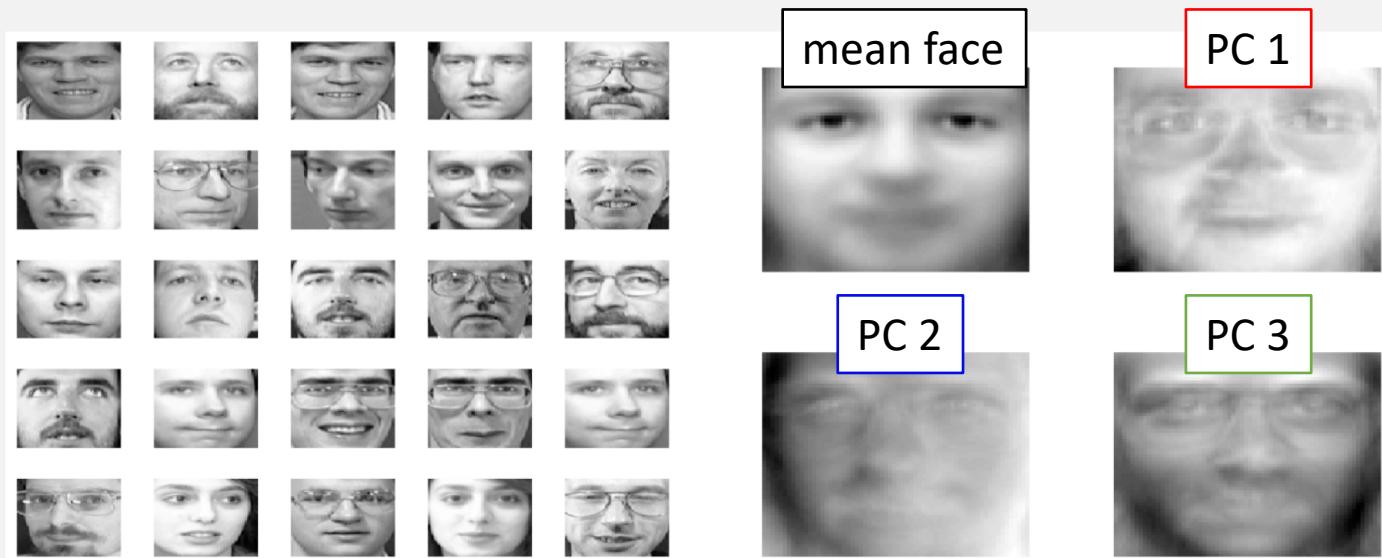
Principle Component Analysis (PCA)

- ▶ Maximise the covariance
- ▶ Linearly independent features

Linear dimensionality reduction



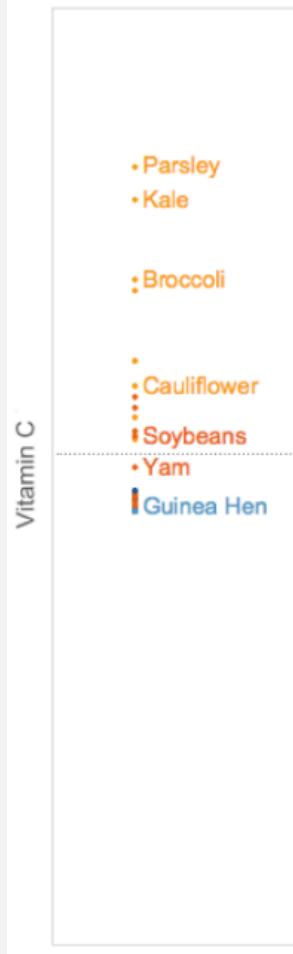
PCA Goal: reduce dimensionality, while preserving the maximum amount of variability



Maximising the variance

We have a dataset of food with the following features ($F=4$)

- ▶ Vitamin C
- ▶ Fat
- ▶ Protein
- ▶ Fiber



Maximising the variance

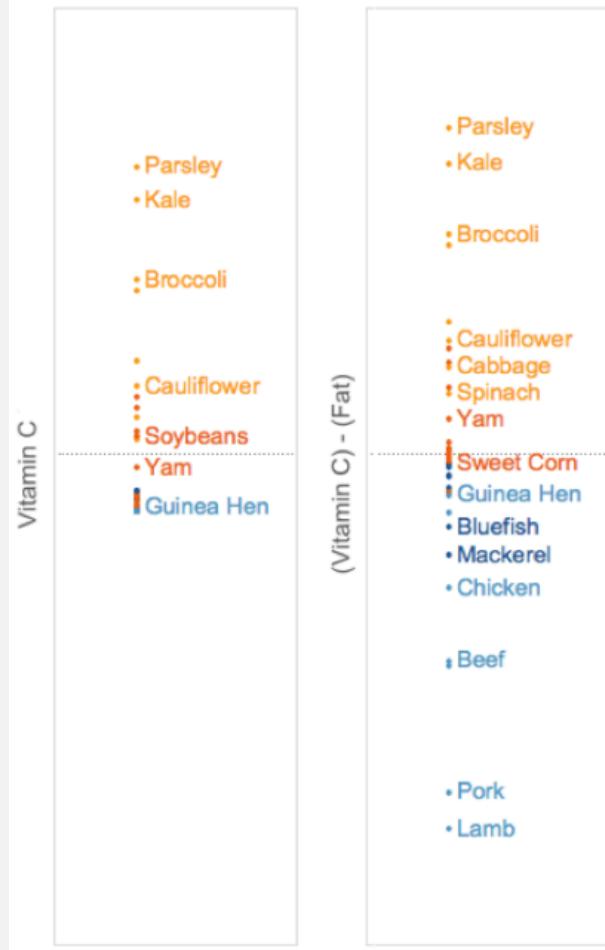
We have a dataset of food with the following features ($F=4$)

- ▶ Vitamin C
- ▶ Fat
- ▶ Protein
- ▶ Fiber

Make a new feature: **VitaminC - Fat**

(taking care to normalise the values to make the comparable)

- ▶ The data spreads out more



Maximising the variance

We have a dataset of food with the following features ($F=4$)

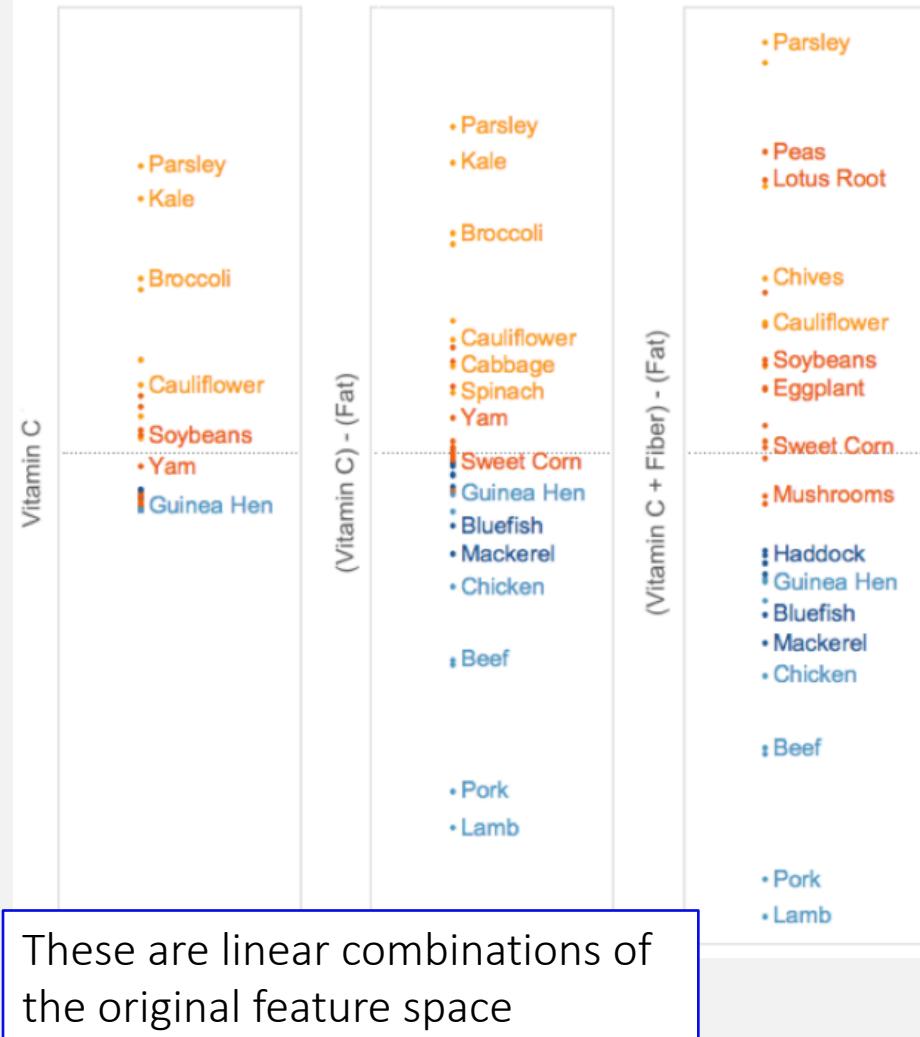
- ▶ Vitamin C
- ▶ Fat
- ▶ Protein
- ▶ Fiber

Make a new feature: $\text{VitaminC} - \text{Fat}$
(taking care to normalise the values to make the comparable)

- ▶ The data spreads out more

Now sum $\text{VitaminC} - \text{Fat} + \text{Fiber}$

- ▶ Even better spread



Principle components and food

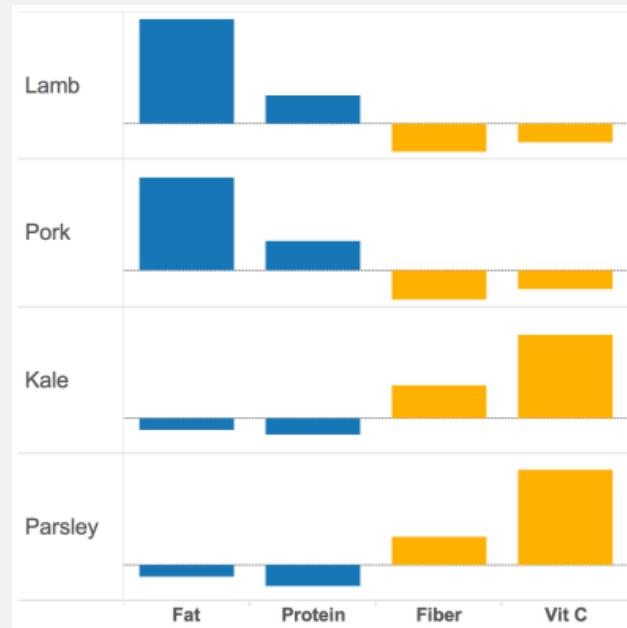
We can see correlations in co-occurrence of some features

- ▶ fat and protein ‘move together’
- ▶ protein and fiber ‘move apart’

Idea:

- ▶ instead of using all 4 features, combine those that are correlated
- ▶ This is the basis of PCA
- ▶ Create new features based on linear combination of the old ones (principle components)
- ▶ Use weights to indicate the importance of each feature’s contribution:

Covariance in features



	PC1	PC2	PC3	PC4
Fat	-0.45	0.66	0.58	0.18
Protein	-0.55	0.21	-0.46	-0.67
Fiber	0.55	0.19	0.43	-0.69
Vitamin C	0.44	0.70	-0.52	0.22

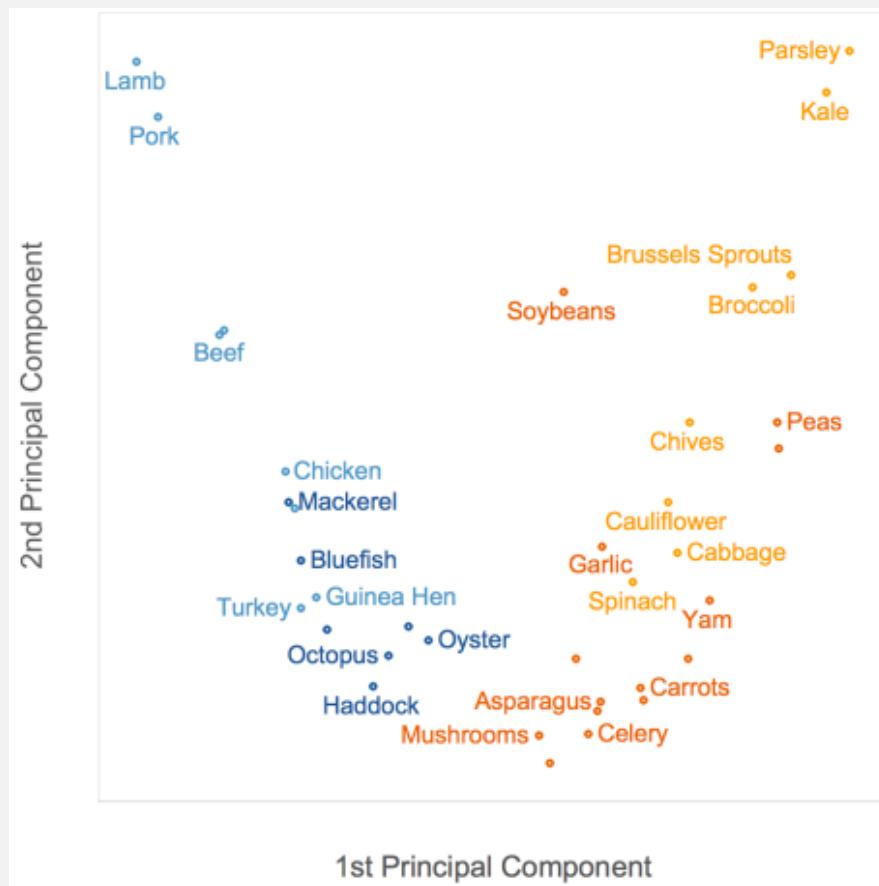
PCA and food

Using our top 2 principle components, we can visualise the data clearly

Calculation:

$$\begin{aligned} \text{PC1} &= 0.55 * \text{fiber} + 0.44 * \text{VitC} \\ &\quad - 0.55 * \text{protein} - 0.45 * \text{fat} \end{aligned}$$

$$\begin{aligned} \text{PC2} &= 0.66 * \text{fat} + 0.7 * \text{VitC} \\ &\quad + 0.21 * \text{protein} + 0.19 * \text{fiber} \end{aligned}$$



	PC1	PC2	PC3	PC4
Fat	-0.45	0.66	0.58	0.18
Protein	-0.55	0.21	-0.46	-0.67
Fiber	0.55	0.19	0.43	-0.69
Vitamin C	0.44	0.70	-0.52	0.22

Principle Component Analysis (PCA)

PCA maximises the covariance in a linearly projected space

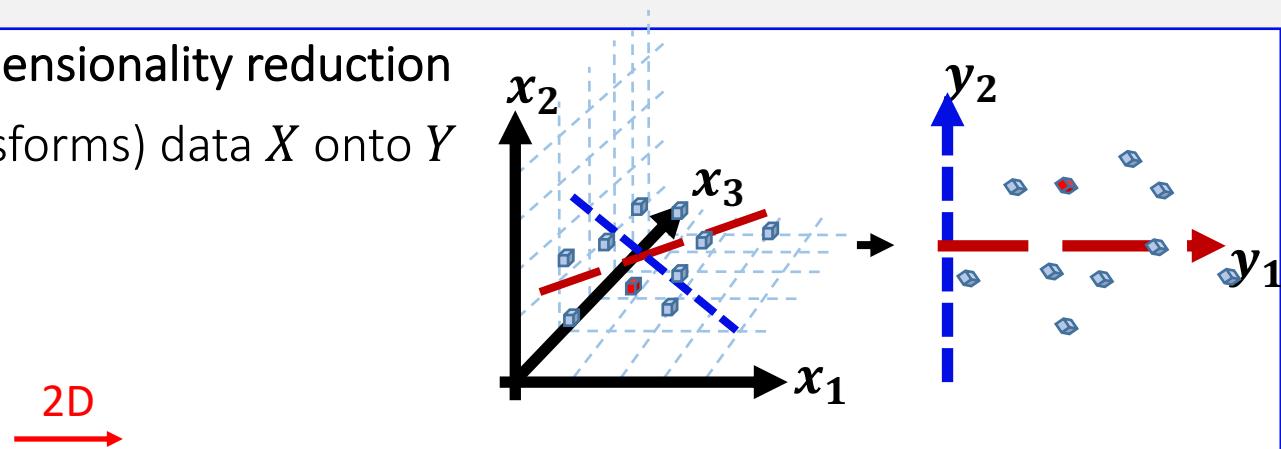
If X represents our data

$$Y = W^T X$$

A linear projection

Linear projection and dimensionality reduction

e.g. matrix W maps (transforms) data X onto Y



If (3 to 2 mapping) $W = \begin{bmatrix} 1 & 0 \\ 0 & 2 \\ 1 & 0 \end{bmatrix}$ and (3x1 feature) data $X = \begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix}$, then

$$Y = W^T X = \begin{bmatrix} 1 & 0 & 1 \\ 0 & 2 & 0 \end{bmatrix} \begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix} = \begin{bmatrix} (1 * 1) + (0 * 2) + (1 * 3) \\ (0 * 1) + (2 * 2) + (0 * 3) \end{bmatrix} = \begin{bmatrix} 4 \\ 4 \end{bmatrix}$$

From 3D to 2D

Linear projection examples

From 3 to 2 dimensions: $W = \begin{bmatrix} 1 & 0 \\ 0 & 2 \\ 1 & 0 \end{bmatrix}$

$$Y = W^T X = \begin{bmatrix} 1 & 0 & 1 \\ 0 & 2 & 0 \end{bmatrix} \begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix} = \begin{bmatrix} (1 * 1) + (0 * 2) + (1 * 3) \\ (0 * 1) + (2 * 2) + (0 * 3) \end{bmatrix} = \begin{bmatrix} 4 \\ 4 \end{bmatrix}$$

From 3 to 1 dimension: $W = \begin{bmatrix} 1 \\ 0 \\ 1 \end{bmatrix}$

$$Y = W^T X = [1 \ 0 \ 1] \begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix} = [(1 * 1) + (0 * 2) + (1 * 3)] = 4$$

From 3 to 3 dimensions: $W = \begin{bmatrix} 1 & 0 & 2 \\ 0 & 2 & 0 \\ 1 & 0 & 2 \end{bmatrix}$

$$Y = W^T X = \begin{bmatrix} 1 & 0 & 1 \\ 0 & 2 & 0 \\ 2 & 0 & 2 \end{bmatrix} \begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix} = \begin{bmatrix} (1 * 1) + (0 * 2) + (1 * 3) \\ (0 * 1) + (2 * 2) + (0 * 3) \\ (2 * 1) + (0 * 2) + (2 * 3) \end{bmatrix} = \begin{bmatrix} 4 \\ 4 \\ 8 \end{bmatrix}$$

Using single 3x1 data point

$$X = \begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix}$$

Principle Component Analysis (PCA)

PCA maximises the covariance in a linearly projected space

If X represents our data

$$Y = W^T X \quad \text{A linear projection}$$
$$\text{cov}(Y) = \text{cov}(W^T X) \quad \text{What we want to maximise}$$

We can write $\text{cov}(Y)$ as a function of $\text{cov}(X)$ directly

$$\text{cov}(Y) = \text{cov}(W^T X) = W^T \text{cov}(X) W$$

With $S = \text{cov}(X)$, we need to find the best W that maximises:

 Square matrix (DxD)

$$W^T S W$$

aside: argmax notation

The argmax notation provides a way to ask

“What is the best value of V that maximises the value of function $f(V)$?”

$$\underset{V}{\operatorname{argmax}}(f(V))$$

Principle Component Analysis (PCA)

PCA maximises the covariance in a linearly projected space

If X represents our data

$$Y = W^T X \quad \text{A linear projection}$$
$$\text{cov}(Y) = \text{cov}(W^T X) \quad \text{What we want to maximise}$$

We can write $\text{cov}(Y)$ as a function of $\text{cov}(X)$ directly

$$\text{cov}(Y) = \text{cov}(W^T X) = W^T \text{cov}(X) W$$

With $S = \text{cov}(X)$, we need to find the best W that maximises:

 Square matrix (DxD)

$$W^T S W$$

That is, find: $\underset{W}{\operatorname{argmax}} (W^T S W)$

Principle Component Analysis (PCA)

PCA maximises the covariance of our data in a linearly projected space

- ▶ Find the best values for W that maximise $W^T SW$:

$$\operatorname{argmax}_W (W^T SW)$$

- ▶ What will happen if you maximise this as it is (think about the 1D case)?

By setting values in W to infinity (∞) we get infinitely big covariance

- ▶ But we wouldn't know what to do with this, so...
- ▶ Add a constraint to ensure W is finite:

$$W^T W = I$$

- ▶ This constraint also ensures **no linear dependencies** in the projected data.
- ▶ Why?...

No linear dependencies?

If the dot product of two (non-zero) vectors is zero, they are linearly independent (orthogonal)

$$\mathbf{x} \cdot \mathbf{y} = \mathbf{x}^T \mathbf{y} = [x_1 \ x_2 \ x_3] \begin{bmatrix} y_1 \\ y_2 \\ y_3 \end{bmatrix} = x_1 y_1 + x_2 y_2 + x_3 y_3 = 0$$

Consider a W with 2 components, where $W^T W = I$:

$$I_2 = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \quad W = [\mathbf{w}_1 \ \mathbf{w}_2], \quad \text{where } \mathbf{w}_i = \begin{bmatrix} w_{i1} \\ w_{i2} \\ \vdots \\ w_{iD} \end{bmatrix}$$

What does the $W^T W$ product look like?

$$W^T W = \begin{bmatrix} \mathbf{w}_1^2 & \mathbf{w}_2^T \mathbf{w}_1 \\ \mathbf{w}_1^T \mathbf{w}_2 & \mathbf{w}_2^2 \end{bmatrix} = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$

Every vector \mathbf{w}_i of W will be linearly independent of all other vectors \mathbf{w}_j

- $\mathbf{w}_i^T \mathbf{w}_j = 0$ (for all $i \neq j$)

Principle Component Analysis (PCA)

PCA maximises the covariance of our data in a linearly projected space

- Find the best values for W that maximise $W^T SW$:

$$\underset{W}{\operatorname{argmax}} (W^T SW) \quad \text{s.t. } W^T W = I$$

Constraint $W^T W = I$ ensures:

- We don't get infinite values for W
- That the projected features are linearly independent

We can use a technique called **eigenanalysis** to get the optimal values of W .

$$[W, \Lambda] = eig(S), \quad \text{where } S = \text{cov}(X)$$

eigenvectors eigenvalues original data

The data can finally be transformed using: $Y = W^T X$

Eigenanalysis (briefly)

$$[\mathbf{W}, \Lambda] = \text{eig}(\mathbf{S}), \quad \text{where } \mathbf{S} = \text{cov}(X)$$

This function finds a matrix of vectors (\mathbf{W}) and corresponding scalars (Λ), such that:

$$\mathbf{SW} = \Lambda\mathbf{W}$$

- multiplying the vectors in \mathbf{W} by covariance matrix \mathbf{S} , is the same as multiplying them by the scalar values Λ
- \mathbf{W} are called the **eigenvectors**
- Λ are called **eigenvalues**

Principle Component Analysis (PCA)

PCA maximises the covariance of our data in a linearly projected space

- Find the best values for W that maximise $W^T S W$:

$$\underset{W}{\operatorname{argmax}} (W^T S W) \quad \text{s.t. } W^T W = I$$

Constraint $W^T W = I$ ensures:

- We don't get infinite values for W
- That the projected features are linearly independent

We can use a technique called **eigenanalysis** to get the optimal values of W .

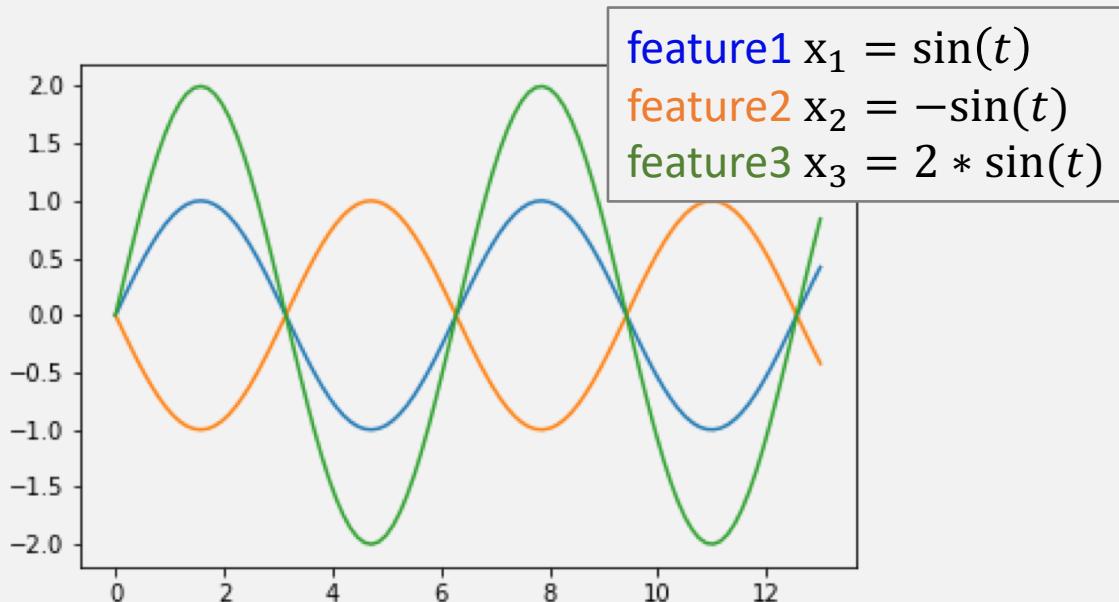
$$[W, \Lambda] = \text{eig}(S), \quad \text{where } S = \text{cov}(X)$$

eigenvectors eigenvalues original data

A diagram showing three red arrows pointing from the text "eigenvectors", "eigenvalues", and "original data" to the equation [W, Λ] = eig(S), where S = cov(X). The word "eigenvectors" has a red arrow pointing to the matrix W. The word "eigenvalues" has a blue arrow pointing to the diagonal matrix Λ. The phrase "original data" has a red arrow pointing to the term cov(X) in the equation.

The data can finally be transformed using: $Y = W^T X$

PCA example

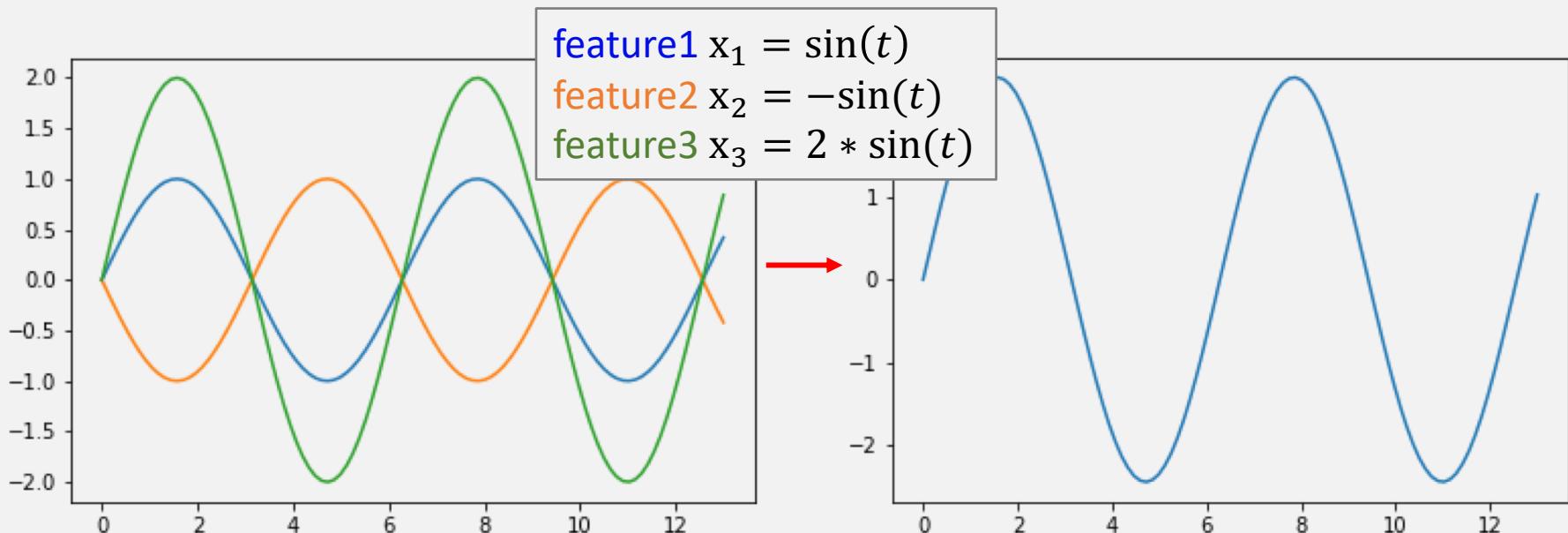


Recall that PCA removes any linear dependencies in the data

Questions:

1. How many dimensions do we need to describe the data above?
2. How many dimensions would PCA need to recover?

PCA example



Recall that PCA removes any linear dependencies in the data

Questions:

1. How many dimensions do we need to describe the data above?
Feature2 is linearly dependent on feature1: $x_2 = -1 * x_1$
Feature3 is linearly dependent on feature2: $x_3 = -2 * x_2 = 2 * x_1$
2. How many dimensions would PCA need to recover?

Applying PCA with Eigenanalysis

Given data X , say in 3D, we want a 2D representation Y that maximises the covariance (using PCA)

1. Centre data around its mean, i.e. $X = X - \bar{X}$
2. Calculate covariance matrix of the data, $S = cov(X)$
3. Calculate eigenvalues and eigenvectors, $[\mathbf{W}, \Lambda] = eig(S)$
4. Order the eigenvalues Λ from large to small → this gives us the components of the corresponding eigenvectors, \mathbf{W}
e.g. for a 2D output, we keep 2 eigenvectors w_i from \mathbf{W} that correspond to the 2 biggest eigenvalues. We can call this matrix \mathbf{W}_2
5. Finally, project the original data into 2D by multiplying \mathbf{W}_2 with X ,
i.e. $Y = \mathbf{W}_2^T X$

```
import numpy as np
X = np.array([[1,2],[.1,.3],[4,3]]) # 2 x 3D data points
S = np.cov(X)
v, U = np.linalg.eig(S)
top = np.flip(np.argsort(v), 0)
U2 = U[:,top[:2]]
Y = U2.T.dot(X)
```

Data compression

Once we have a feature reduction map \mathbf{W} :

$$Y_2 = \mathbf{W}_2^T X$$

e.g. \mathbf{W}_2^T X Y_2

$$= \begin{bmatrix} .7 & -.7 & .2 \\ -.14 & .14 & .98 \end{bmatrix} \begin{bmatrix} 1 \\ -1 \\ .5 \end{bmatrix} \approx \begin{bmatrix} 1.44 \\ 0.23 \end{bmatrix}$$

We can reverse the operation to recreate our original data:

$$\tilde{X} = \mathbf{W}_2 Y_2$$

\mathbf{W}_2 \tilde{X}

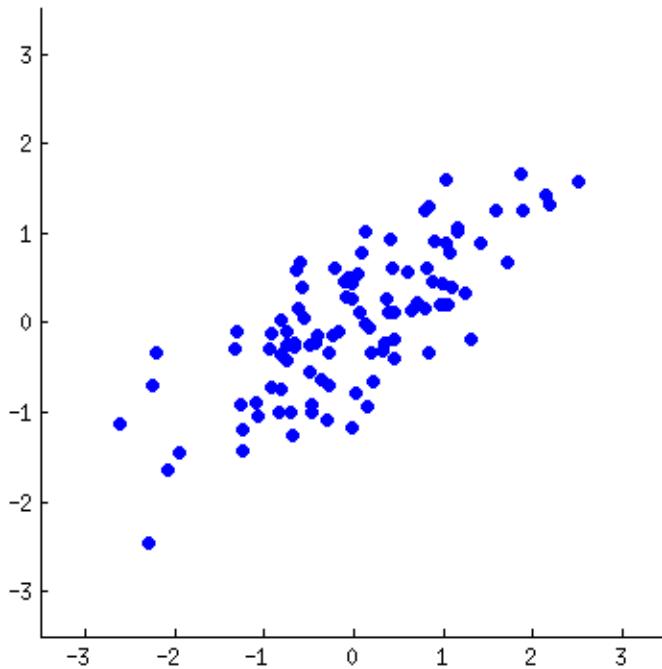
$$= \begin{bmatrix} .7 & -.14 \\ -.7 & .14 \\ .2 & .98 \end{bmatrix} \begin{bmatrix} 1.44 \\ 0.23 \end{bmatrix} \approx \begin{bmatrix} 1 \\ -1 \\ .5 \end{bmatrix}$$

Some data will be lost (unless we retain all components), but PCA ensures the most descriptive data remains.

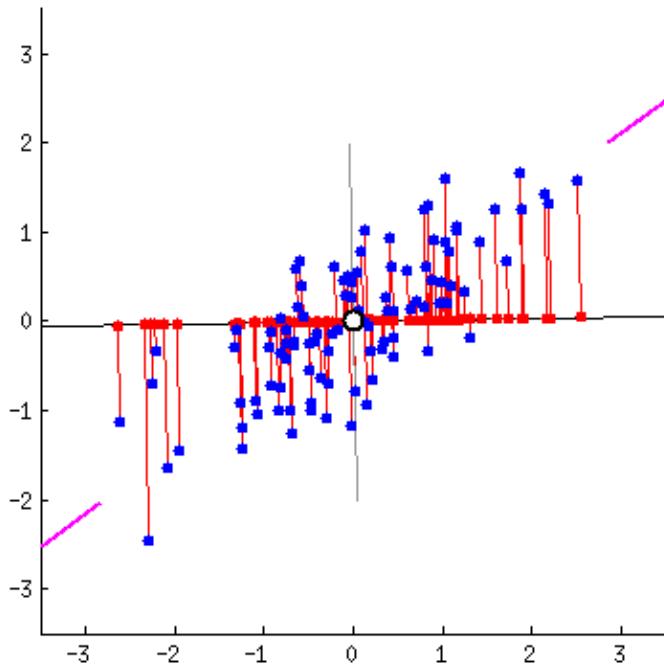
Same as above, but using a 1D mapping
 $\tilde{X} = \mathbf{W}_1 Y_1$

$$\begin{bmatrix} .7 \\ -.7 \\ .2 \end{bmatrix} [1.44] \approx \begin{bmatrix} 1 \\ -1 \\ .3 \end{bmatrix}$$

PCA: maximise variance, minimise error



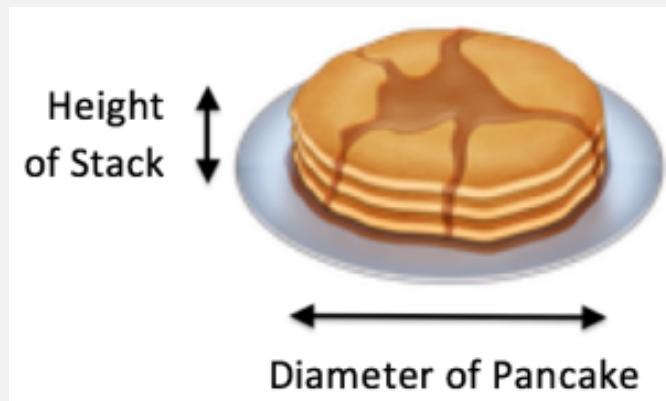
PCA: maximise variance, minimise error



Limitation: the pancake problem

Maximising variance isn't always what we need

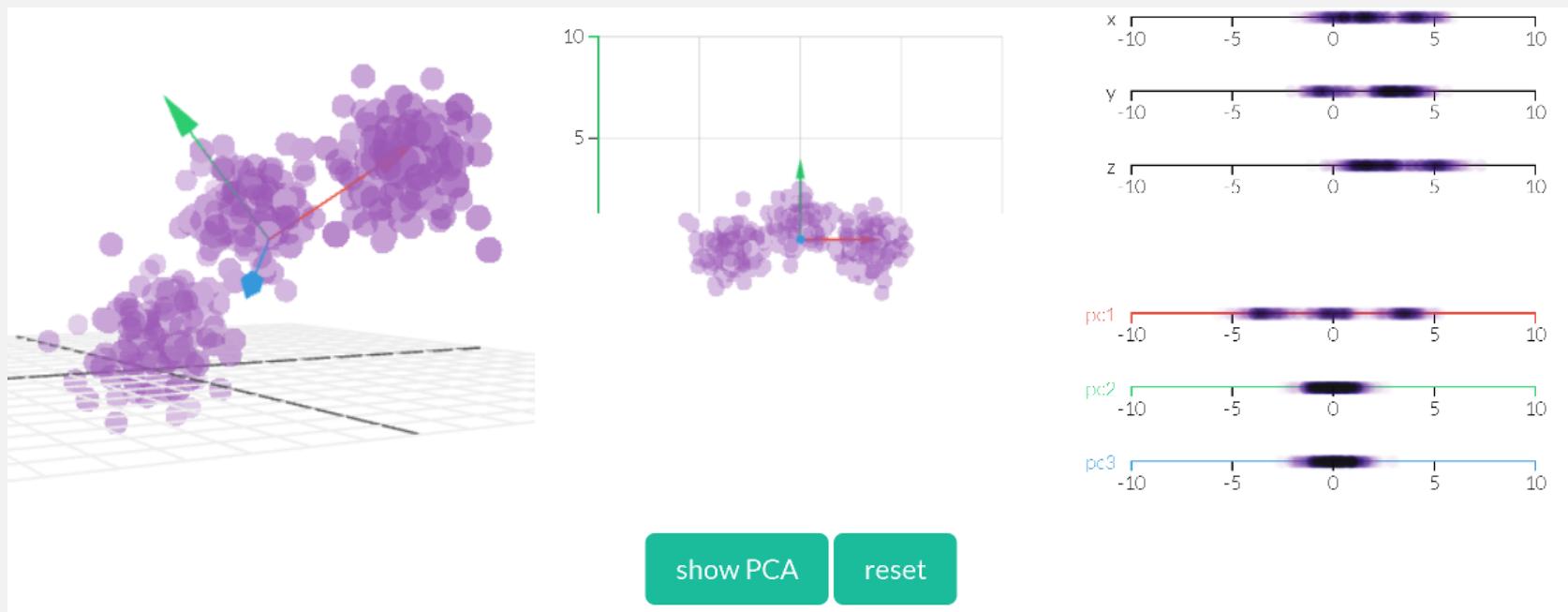
- ▶ If we are interested in the number of pancakes, then height is a more useful dimension than diameter
- ▶ But PCA will select diameter if it has more variance (i.e. for a small number of pancakes)



Visualisations

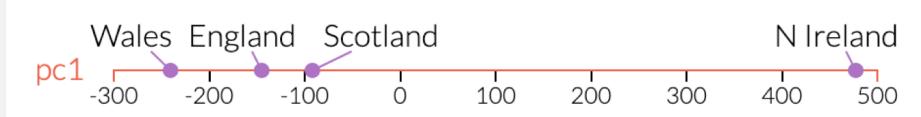
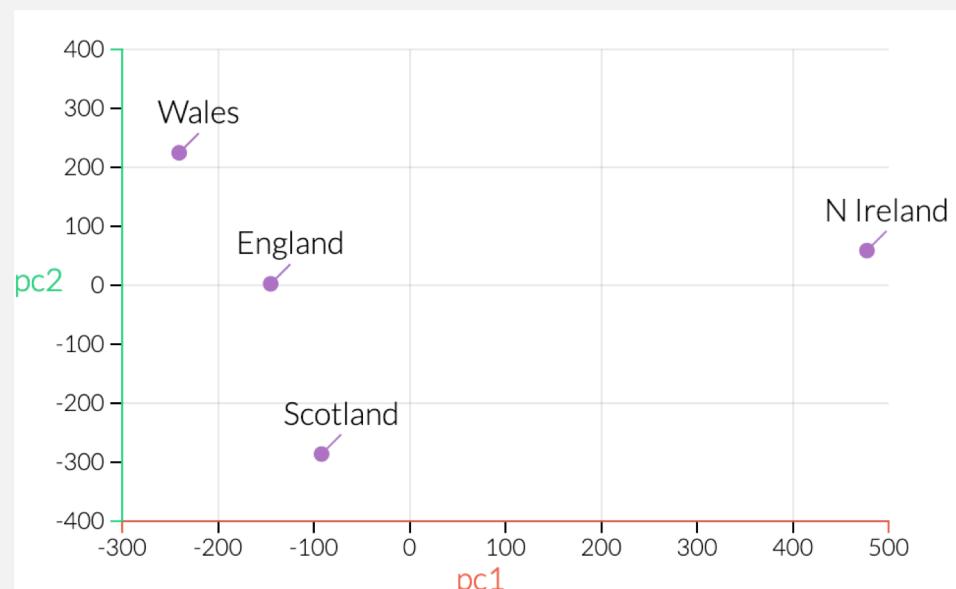
Some nice interactive examples:

<http://setosa.io/ev/principal-component-analysis/>



Eating in the UK (a 17D example)

	England	N Ireland	Scotland	Wales
Alcoholic drinks	375	135	458	475
Beverages	57	47	53	73
Carcase meat	245	267	242	227
Cereals	1472	1494	1462	1582
Cheese	105	66	103	103
Confectionery	54	41	62	64
Fats and oils	193	209	184	235
Fish	147	93	122	160
Fresh fruit	1102	674	957	1137
Fresh potatoes	720	1033	566	874
Fresh Veg	253	143	171	265
Other meat	685	586	750	803
Other Veg	488	355	418	570
Processed potatoes	198	187	220	203
Processed Veg	360	334	337	365
Soft drinks	1374	1506	1572	1256
Sugars	156	139	147	175



Interactive examples:

<http://setosa.io/ev/principal-component-analysis/>