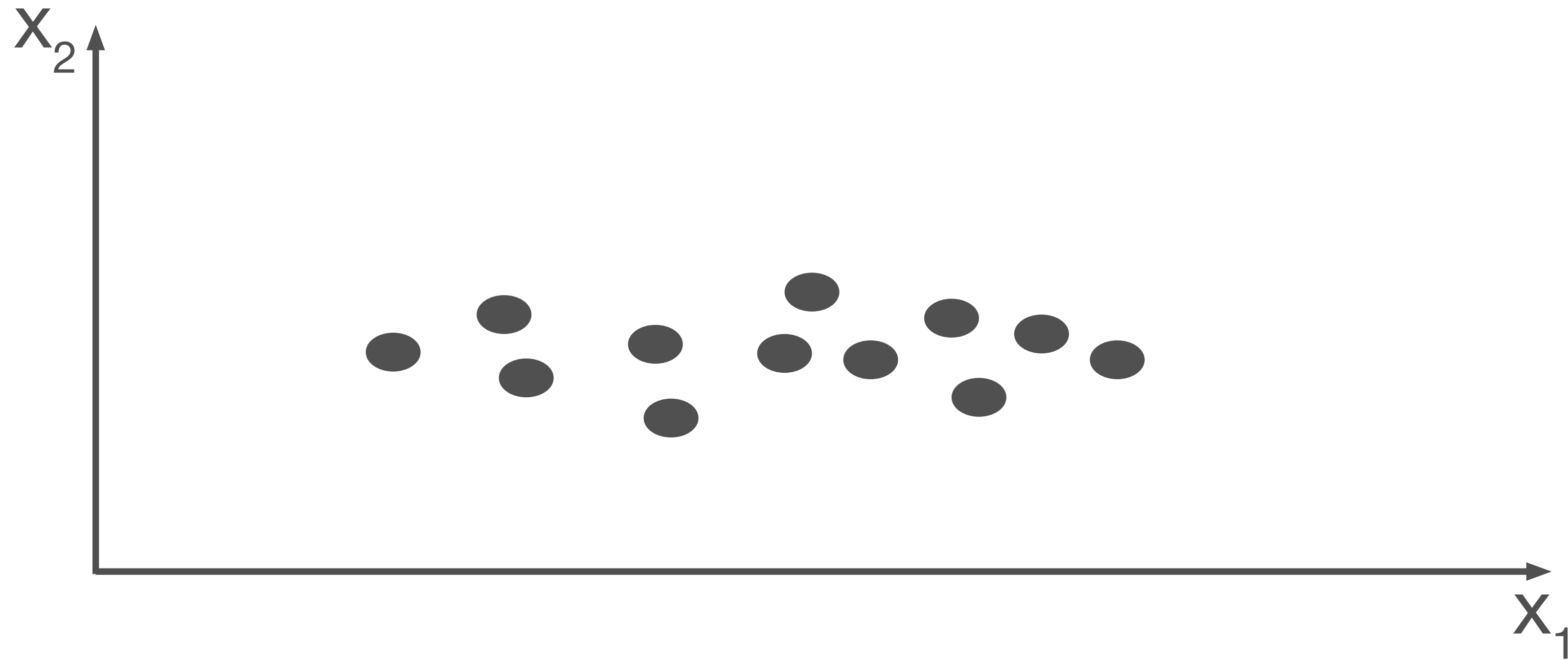# Unsupervised learning

- no target variable

- multidimensional data

- usually a part of a larger task
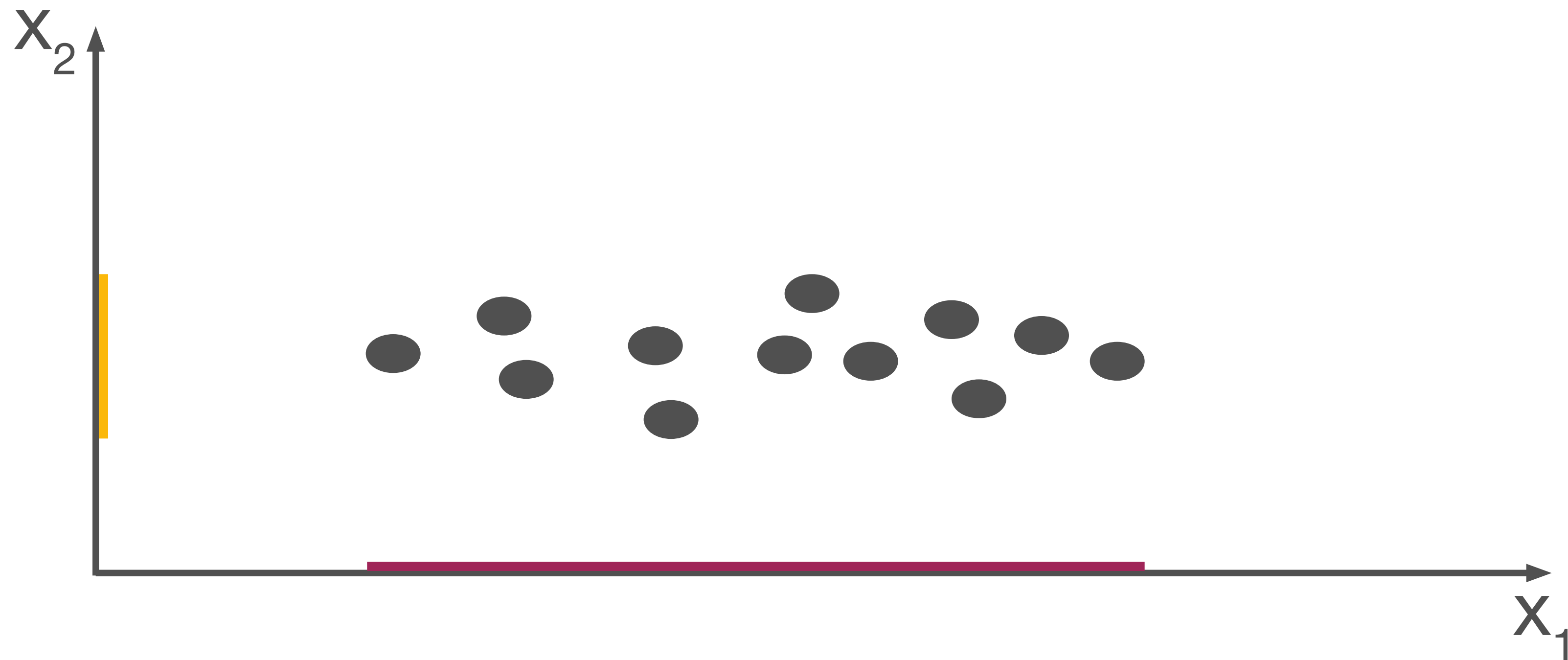
# How to approach

- What is the goal/question?
  - Can we reduce dimensions while preserving `information`?
  - Are there clusters in the data?
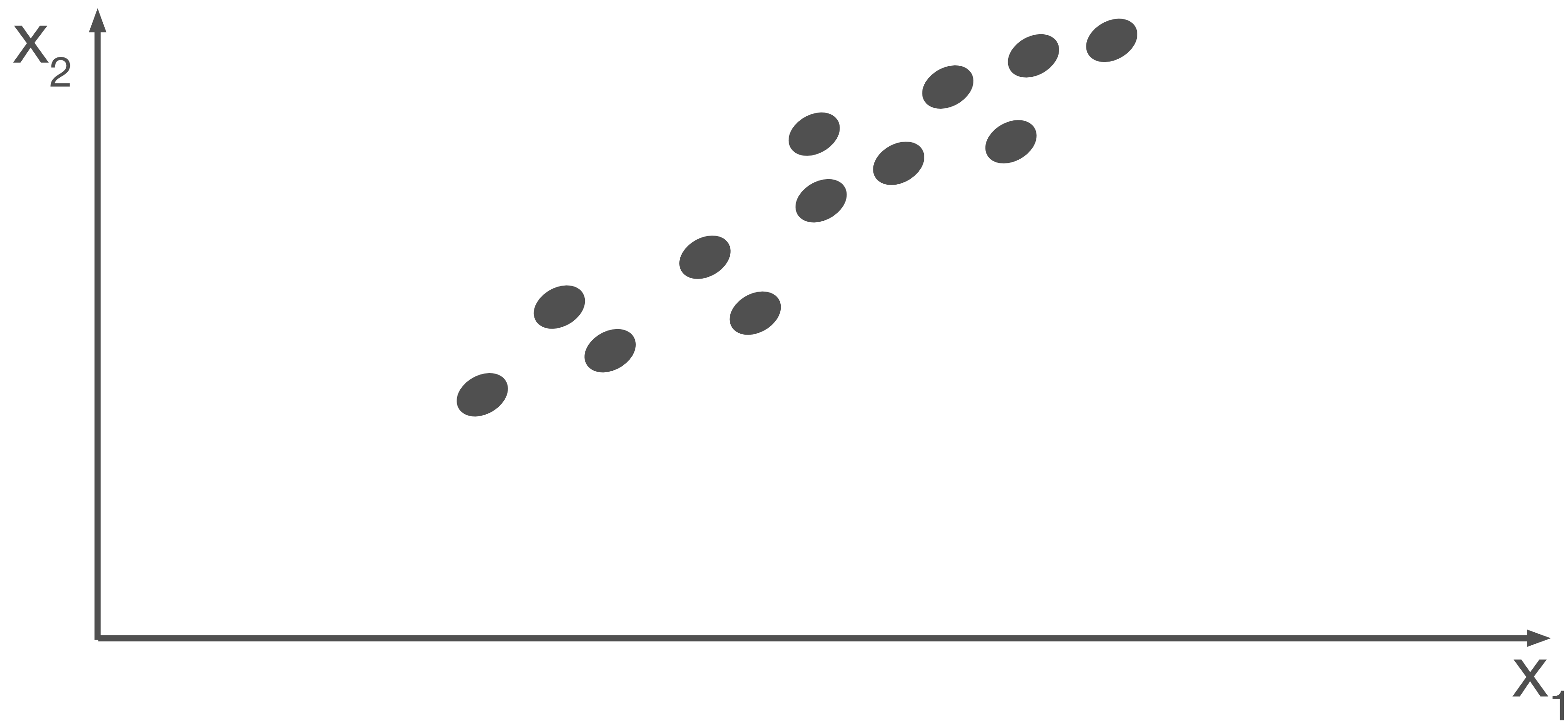
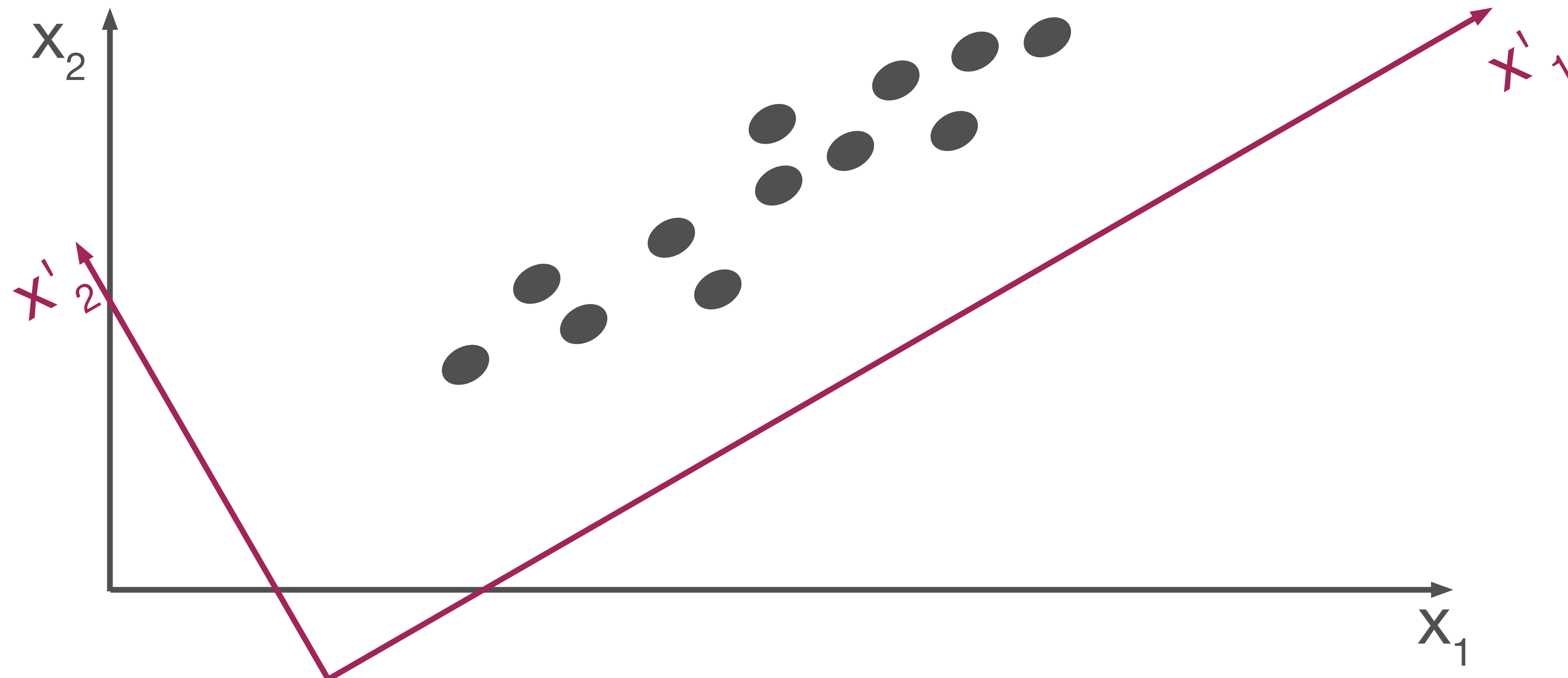- Select appropriate method for the question.

# Goal:

- Can we reduce dimensions, but preserve `information`?

# PCA: intuition

- change axes
- preserve variance

# Variance as a measure of `information`

$$Cov(A, B) = \frac{1}{n-1} \sum_{i=1}^{n} (A_i - \bar{A})(B_i - \bar{B})$$

$$Var(A) = Cov(A, A)$$

# Matrix form

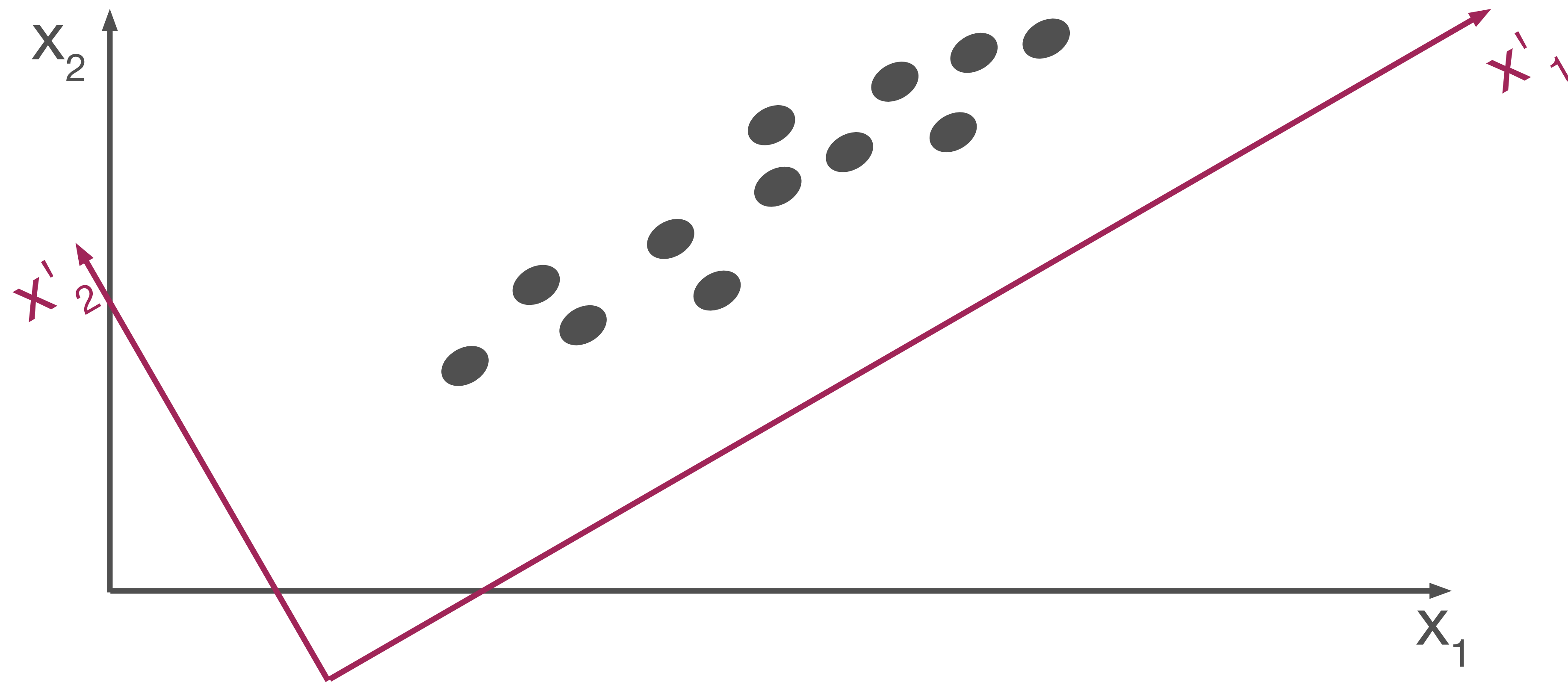$$X = \begin{pmatrix} a_1 & b_1 & c_1 \\ a_2 & b_2 & c_2 \\ a_3 & b_3 & c_3 \\ a_4 & b_4 & c_4 \end{pmatrix}$$

$$CovMatrix = \frac{1}{n-1}(X - \bar{X})^t(X - \bar{X}) = \begin{pmatrix} cov_{11} & cov_{12} & cov_{13} \\ cov_{21} & cov_{22} & cov_{23} \\ cov_{31} & cov_{32} & cov_{33} \end{pmatrix}$$
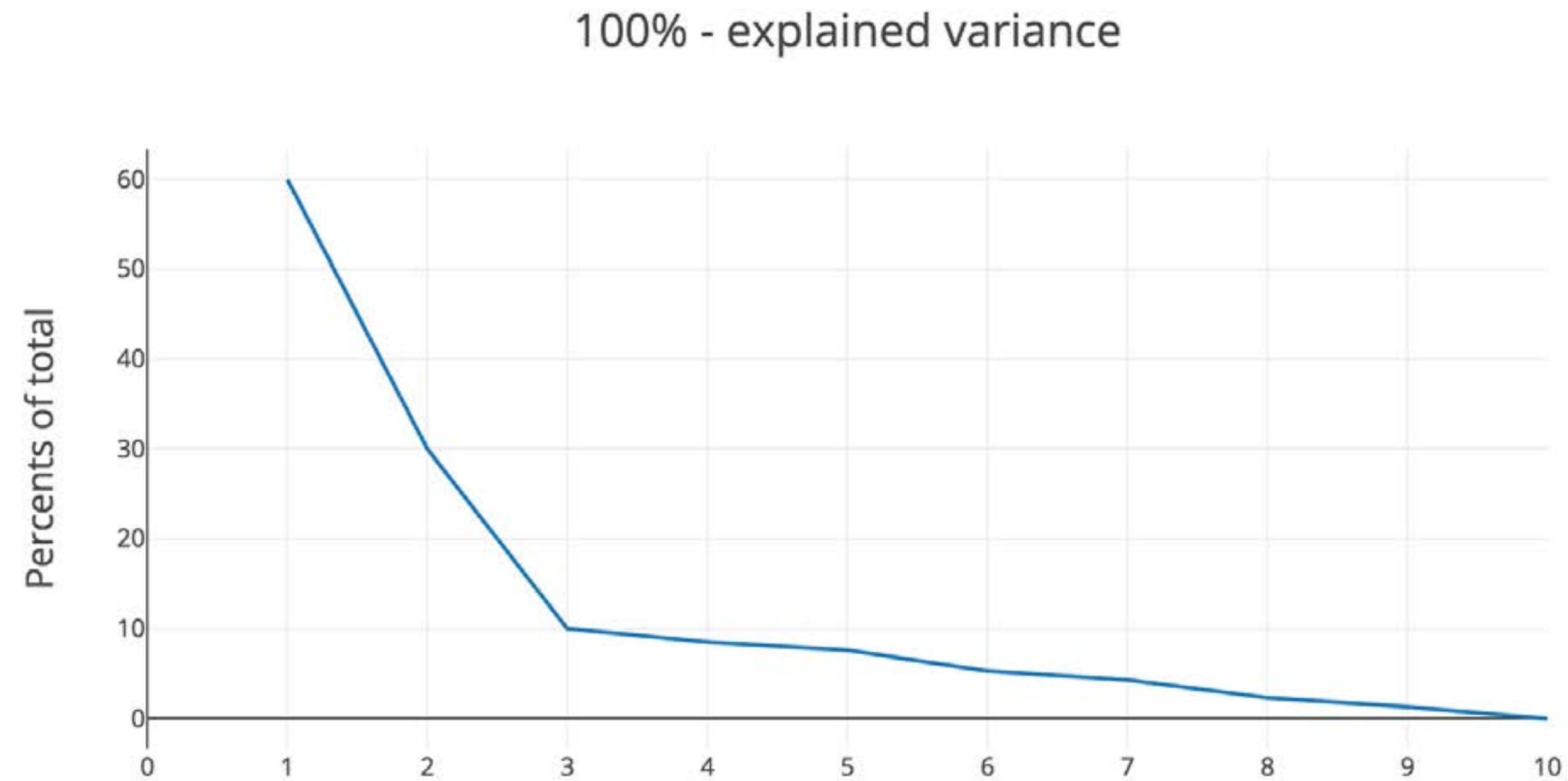
# Eigenvalue decomposition

$$\begin{pmatrix} cov_{11} & cov_{12} & cov_{13} \\ cov_{21} & cov_{22} & cov_{23} \\ cov_{31} & cov_{32} & cov_{33} \end{pmatrix} =$$

$$\begin{pmatrix} f_1^1 & f_1^2 & f_1^3 \\ f_2^1 & f_2^2 & f_2^3 \\ f_3^1 & f_3^2 & f_3^3 \end{pmatrix} \begin{pmatrix} \lambda_1 & 0 & 0 \\ 0 & \lambda_2 & 0 \\ 0 & 0 & \lambda_3 \end{pmatrix} \begin{pmatrix} f_1^1 & f_2^1 & f_3^1 \\ f_1^2 & f_2^2 & f_3^2 \\ f_1^3 & f_2^3 & f_3^3 \end{pmatrix}$$

Eigenvalue

Eigenvector

- PCA for reduction
  - Get principal components
  - Take sufficient number of them



100% - explained variance

- Note: apply after scaling

- PCA to reduce dimensions
  - get new axes
  - project on a subset
- imply that preserving `information` means preserving `variance`
  - might want to preserve something else
- not the only way to reduce dimensions