

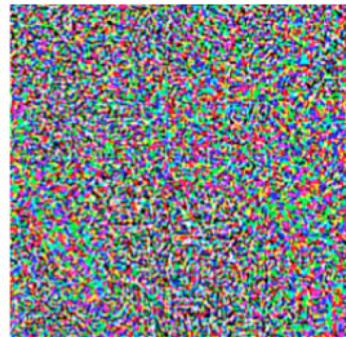
# Lipschitz constant of a convolutional layer in neural network

It was observed, that small perturbation in Neural Network input could lead to significant errors, i.e. misclassifications. Picture below from the [article](#).



**Bagle**

+ 0.001 ×



=



**piano**



**stop sign**

+ 0.001 ×



=



**teddy bear**

Lipschitz constant bounds the magnitude of the output of a function, so it cannot change drastically with a slight change in the input

$$\|NN(image) - NN(image + \varepsilon)\| \leq L\|\varepsilon\|$$

In this notebook we will try to estimate Lipschitz constant of some convolutional layer of a Neural Network.

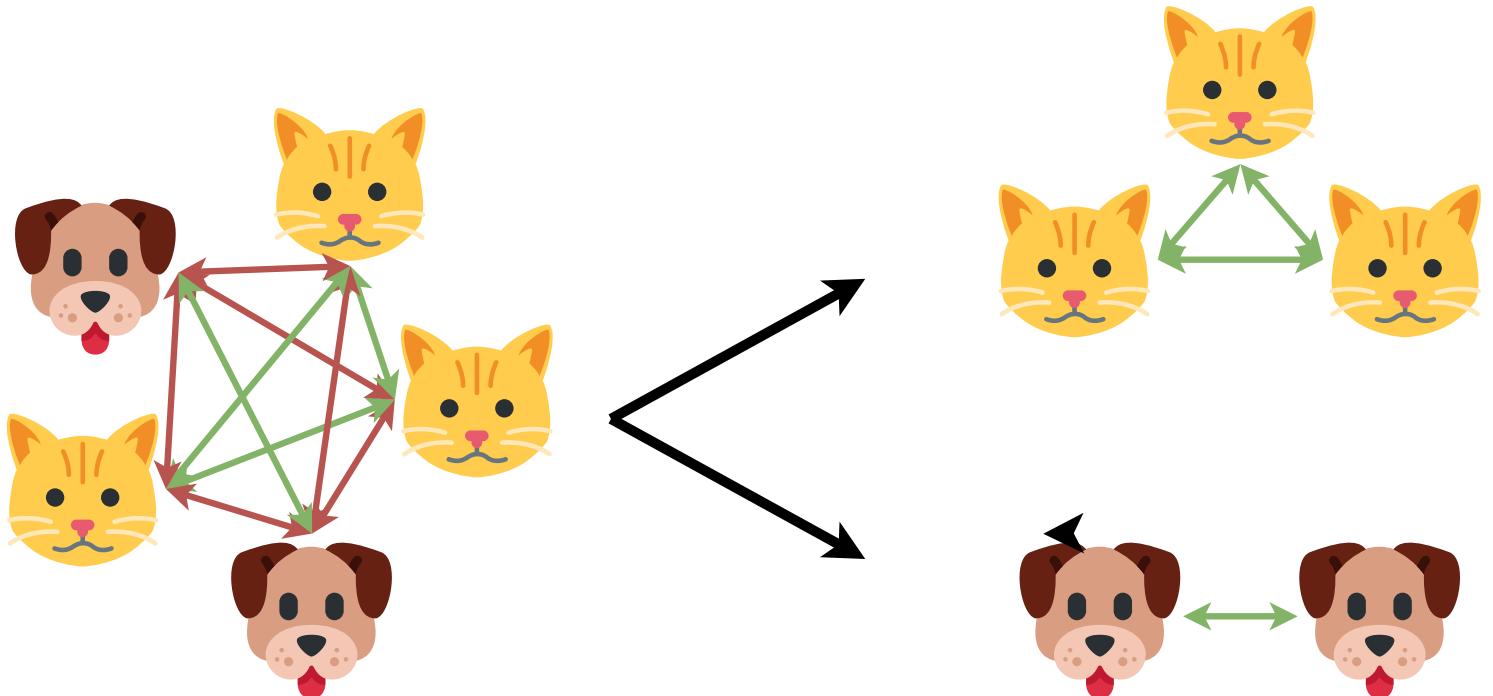
Power method for the singular value estimation of matrix  $W$ :

$$x_{k+1} = \frac{W^\top W x_k}{\|W^\top W x_k\|}$$

$$\sigma_{k+1} = \frac{\|W x_k\|}{\|x_k\|}$$

# Two way partitioning problem

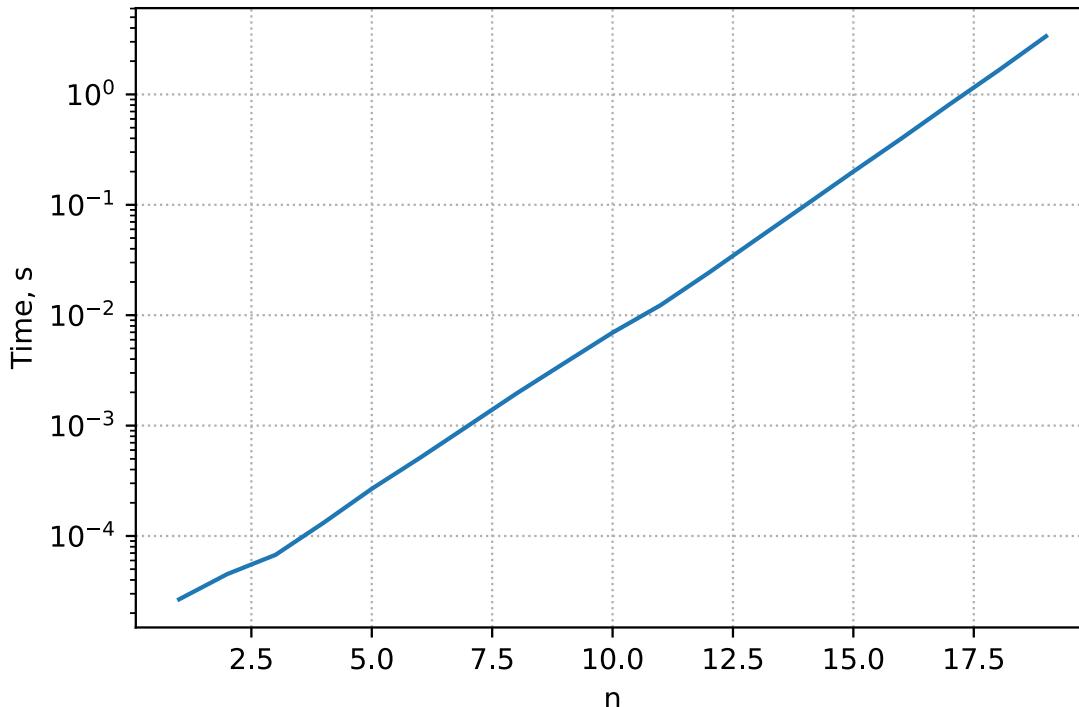
## Intuition



Suppose, we have a set of  $n$  objects, which are needed to be splitted into two groups. Moreover, we have information about the preferences of all possible pairs of objects to be in the same group. this information could be presented in the matrix form:  $W \in \mathbb{R}^{n \times n}$ , where  $\{w_{ij}\}$  is the cost of having  $i$ -th and  $j$ -th object in the same partitions. It is easy to see, that the total number of partitions is finite and equals to  $2^n$ . So this problem can in principle be solved by simply checking the objective value of each feasible point. Since the number of feasible points grows exponentially, however, this is possible only for small problems (say, with  $n \leq 30$ ). In general (and for  $n$  larger than, say, 50) the problem is very difficult to solve.

For example, bruteforce solution on MacBook Air with M1 processor without any explicit parallelization will take more, than a universe lifetime for  $n = 62$ .

Average time for brutforce solution. 3 runs per n



Despite the hardness of the problems, there are several ways to approach it.

# Problem

We consider the (nonconvex) problem

$$\begin{aligned} & \min_{x \in \mathbb{R}^n} x^\top W x, \\ & \text{s.t. } x_i^2 = 1, \quad i = 1, \dots, n \end{aligned}$$

where  $W \in \mathbb{R}^n$  is the symmetric matrix. The constraints restrict the values of  $x_i$  to 1 or  $-1$ , so the problem is equivalent to finding the vector with components  $\pm 1$  that minimizes  $x^\top W x$ . The feasible set here is finite (it contains  $2^n$  points), thus, is non-convex.

The objective is the total cost, over all pairs of elements, and the problem is to find the partition with least total cost.

# Simple lower bound with duality

We now derive the dual function for this problem. The Lagrangian is

$$L(x, \nu) = x^\top W x + \sum_{i=1}^n \nu_i (x_i^2 - 1) = x^\top (W + \text{diag}(\nu)) x - \mathbf{1}^\top \nu.$$

We obtain the Lagrange dual function by minimizing over  $x$ :

$$\begin{aligned} g(\nu) &= \inf_{x \in \mathbb{R}^n} x^\top (W + \text{diag}(\nu)) x - \mathbf{1}^\top \nu = \\ &= \begin{cases} \mathbf{1}^\top \nu, & W + \text{diag}(\nu) \succeq 0 \\ -\infty, & \text{otherwise} \end{cases} \end{aligned}$$

so

This dual function provides lower bounds on the optimal value of the difficult problem. For example, we can take any specific value of the dual variable

$$\nu = -\lambda_{\min}(W)\mathbf{1},$$

This yields the bound on the optimal value  $p^*$ :

$$p^* \geq g(\nu) \geq -\mathbf{1}^\top \nu = n\lambda_{\min}(W)$$

**Question** Can you obtain the same lower bound without knowledge of duality, but using the idea of eigenvalues?

## Code

 Open in Colab

## References

- Convex Optimization book by Stephen Boyd and Lieven Vandenberghe.

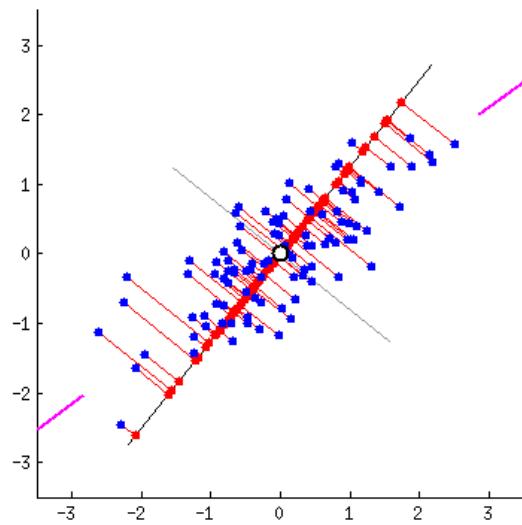
## Eigenfaces

## PCA recap

# Intuition

Imagine, that you have a dataset of points. Your goal is to choose orthogonal axes, that describe your data the most informative way. To be precise, we choose first axis in such a way, that maximize the variance (expressiveness) of the projected data. All the following axes have to be orthogonal to the previously chosen ones, while satisfy largest possible variance of the projections.

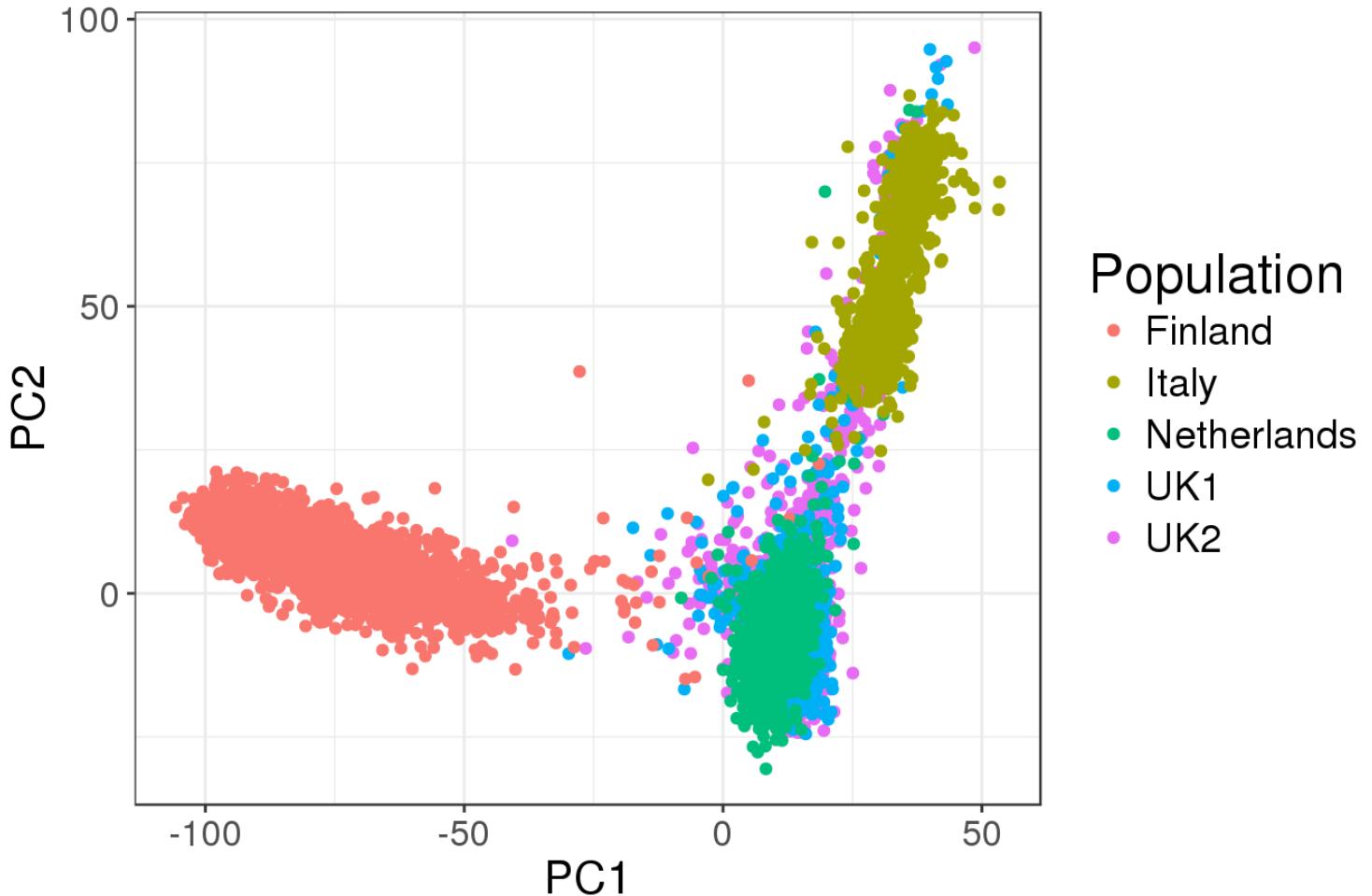
Let's take a look at the simple 2d data. We have a set of blue points on the plane. We can easily see that the projections on the first axis (red dots) have maximum variance at the final position of the animation. The second (and the last) axis should be orthogonal to the previous one.



[source](#)

This idea could be used in a variety of ways. For example, it might happen, that projection of complex data on the principal plane (only 2 components) bring you enough intuition for clustering. The picture below plots projection of the labeled dataset onto the first two principal components (PCs), we can clearly see, that only two vectors (these PCs) would be enough to differ Finnish people from Italian in particular dataset (celiac disease (Dubois et al. 2010))

# Scores of PCA



[source](#)

## Problem

The first component should be defined in order to maximize variance. Suppose, we've already normalized the data, i.e.  $\sum_i a_i = 0$ , then sample variance will become the sum of all squared projections of data points to our vector  $\mathbf{w}_{(1)}$ , which implies the following optimization problem:

$$\mathbf{w}_{(1)} = \arg \max_{\|\mathbf{w}\|=1} \left\{ \sum_i \left( \mathbf{a}_{(i)}^\top \cdot \mathbf{w} \right)^2 \right\}$$

or

$$\mathbf{w}_{(1)} = \arg \max_{\|\mathbf{w}\|=1} \{ \|\mathbf{A}\mathbf{w}\|^2 \} = \arg \max_{\|\mathbf{w}\|=1} \{ \mathbf{w}^\top \mathbf{A}^\top \mathbf{A} \mathbf{w} \}$$

since we are looking for the unit vector, we can reformulate the problem:

$$\mathbf{w}_{(1)} = \arg \max \left\{ \frac{\mathbf{w}^\top \mathbf{A}^\top \mathbf{A} \mathbf{w}}{\mathbf{w}^\top \mathbf{w}} \right\}$$

It is known, that for positive semidefinite matrix  $A^\top A$  such vector is nothing else, but eigenvector of  $A^\top A$ , which corresponds to the largest eigenvalue. The following components will give you the same results (eigenvectors).

So, we can conclude, that the following mapping:

$$\Pi_{n \times k} = A_{n \times d} \cdot W_{d \times k}$$

describes the projection of data onto the  $k$  principal components, where  $W$  contains first (by the size of eigenvalues)  $k$  eigenvectors of  $A^\top A$ .

Now we'll briefly derive how SVD decomposition could lead us to the PCA.

Firstly, we write down SVD decomposition of our matrix:

$$A = U \Sigma W^\top$$

and to its transpose:

$$\begin{aligned} A^\top &= (U \Sigma W^\top)^\top \\ &= (W^\top)^\top \Sigma^\top U^\top \\ &= W \Sigma^\top U^\top \\ &= W \Sigma U^\top \end{aligned}$$

Then, consider matrix  $AA^\top$ :

$$\begin{aligned} A^\top A &= (W \Sigma U^\top)(U \Sigma V^\top) \\ &= W \Sigma I \Sigma W^\top \\ &= W \Sigma \Sigma W^\top \\ &= W \Sigma^2 W^\top \end{aligned}$$

Which corresponds to the eigendecomposition of matrix  $A^\top A$ , where  $W$  stands for the matrix of eigenvectors of  $A^\top A$ , while  $\Sigma^2$  contains eigenvalues of  $A^\top A$ .

At the end:

$$\begin{aligned} \Pi &= A \cdot W = \\ &= U \Sigma W^\top W = U \Sigma \end{aligned}$$

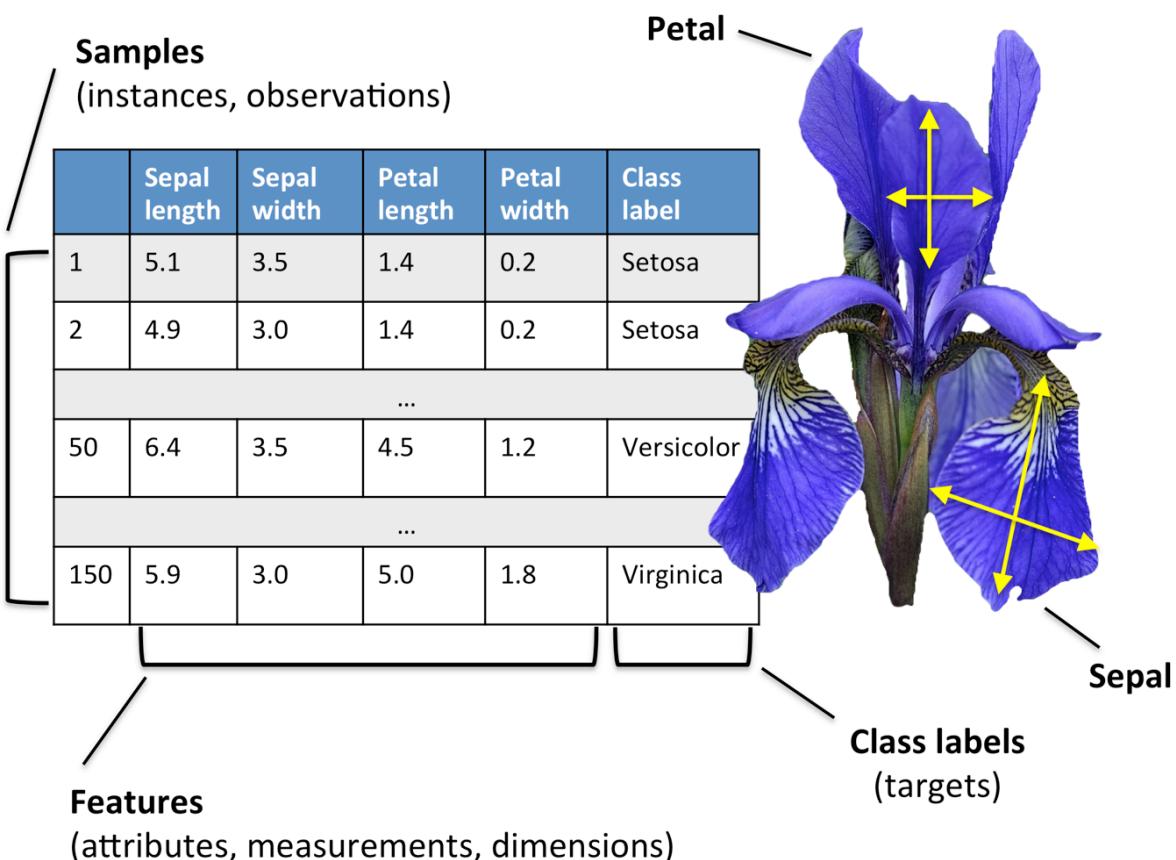
The latter formula provide us with easy way to compute PCA via SVD with any number of principal components:

$$\Pi_r = U_r \Sigma_r$$

## Examples

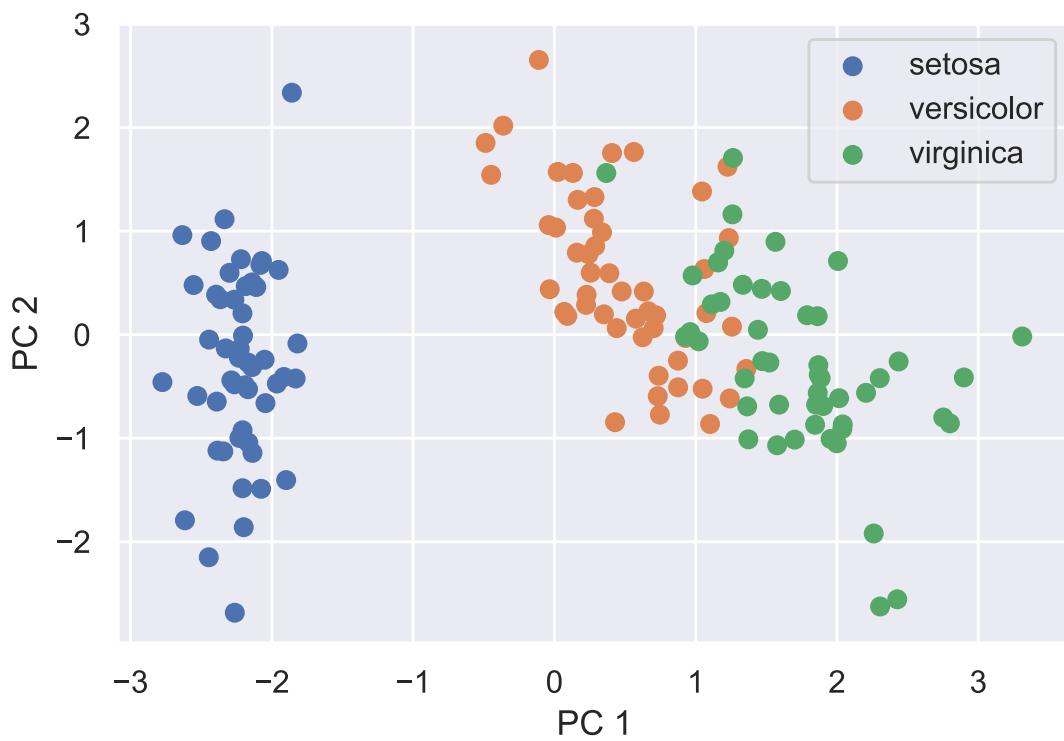
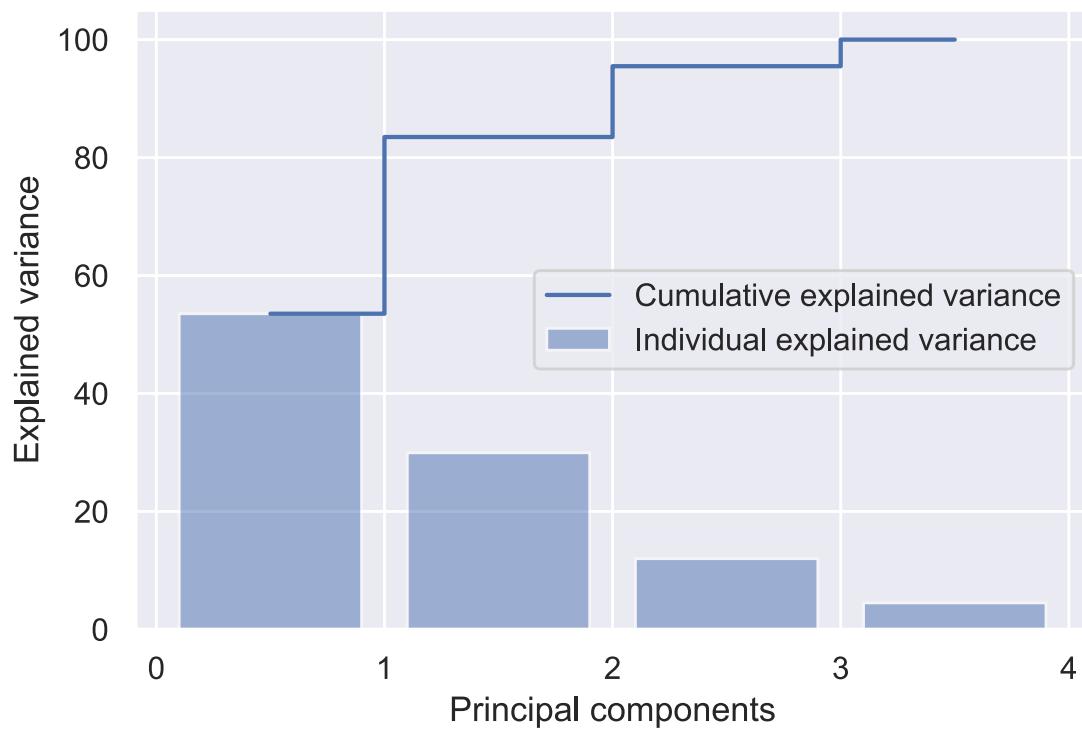
### Iris dataset

Consider the classical Iris dataset



[source](#)

We have the dataset matrix  $A \in \mathbb{R}^{150 \times 4}$



## Code

Open in Colab

## Related materials

- [Wikipedia](#)

- [Blog post](#)
- [Blog post](#)