# Outline

# Eigenvalue Decomposition (EVD)

Eigenvalue decomposition places significant role in PCA. PCs are eigenvectors of $X^\top X$ and PCA properties are derived from those of eigenvectors and eigenvalues.

- For any $p \times p$ *symmetric* matrix $S$ (think for example $X^\top X$), there exists $p$ eigenvectors $v_1, \ldots, v_p$ that are pairwise orthogonal and $p$ associated eigenvalues $\lambda_1, \ldots, \lambda_p$ which satisfy the eigenvalue equation $Sv_i = \lambda_i v_i \; \forall i$.
- $S$ can be written as $S = V\Lambda V^\top$ where
  - $V = [v_1, \ldots, v_p]$ is a $p \times p$ orthogonal matrix
  - $\Lambda = diag\{\lambda_1, \ldots, \lambda_p\}$
  - and if $S_{ij} \in \mathbb{R} \; \forall i,j$, $\lambda_i \in \mathbb{R} \; \forall i$
- The relevant R-command is `eigen`. Look at `?eigen` to get help on the command.

# Singular Value Decomposition (SVD)

The SVD of a matrix $X$ is an equally useful matrix factorisation that is related to the EVD.

- ▶ Though the EVD does not exist for $\mathbb{R}^{n \times p}$ matrices if $p \neq n$, SVDs *always* exists.
- ▶ $X$ can be written as $X = UDV^\top$ where
    - ▶ $U$ is an $n \times n$ matrix with orthogonal columns.
    - ▶ $D$ is a $n \times p$ matrix with decreasing non-negative elements on the diagonal (the singular values) and zero off-diagonal elements.
    - ▶ $V$ is a $p \times p$ matrix with orthogonal columns.

    The relevant R-command is `svd`.
- ▶ SVD can be computed using very fast and numerically stable algorithms.

# Some Properties of the SVD

▶ Let $X = UDV^\top$ be again the SVD of the $n \times p$ matrix $X$.

▶ Note that

$$X^\top X = (UDV^\top)^\top(UDV^\top) = VD^\top U^\top UDV^\top = VD^\top DV^\top,$$

using orthogonality $(U^\top U = I_n)$ of $U$.

▶ The eigenvalues of $S = X^\top X$ are thus the squares of the singular values of $X$ and the columns of the orthogonal matrix $V$ are the eigenvectors of $S$.

▶ We also have

$$XX^\top = (UDV^\top)(UDV^\top)^\top = UDV^\top VD^\top U^\top = UDD^\top U^\top,$$

using orthogonality $(V^\top V = I_p)$ of $V$.

▶ Consider the following optimization problem:

$$\min_{\tilde{X}} \|\tilde{X} - X\|^2 \qquad \text{s.t. } \tilde{X} \text{ has maximum rank } r < n, p.$$

This problem can be solved by keeping only the $r$ largest singular values of $X$, zeroing out the smaller singular values in the SVD.