# Principal Component Algorithm:

**Setup:** $X$ is a dataset : $n \times p$ matrix.

**Notation:**

$E_{1:K}$ is the matrix created from $E$ by discarding the last $n-K$ columns.

$$(E) \quad (E_{1:K} \boxed{\phantom{/}})$$
$\uparrow$ discarded.

**Step #1:**

Center the matrix $X$ by subtracting the column means.

Assumption: From now on, $X$ is centered.

**Step #2:**

Compute the sample covariance matrix $\Omega = \frac{1}{n} X^t X$.

**Step #3:**

Compute the eigenvectors $\{e_1, e_2, \dots, e_p\}$
and eigenvalues $\{\lambda_1, \lambda_2, \dots, \lambda_p\}$ of $\Omega$.

**Return:**

each column is an eigenvector

The matrix of eigenvectors $E = \begin{pmatrix} | & | & & | \\ e_1 & e_2 & \cdots & e_p \\ | & | & & | \end{pmatrix}$

The vector of eigenvalues $\Gamma = (\lambda_1, \lambda_2 \cdots \lambda_p)$

**Properties:**

① The projection of $X$ onto the subspace spanned by $\{e_1, e_2, \dots, e_p\}$ preserves the **maximum** variance (out of all $K$-dim subspaces).

② The matrix product $X E_{1:K}$ gives the coordinates of $X$ in the reduced basis $\{e_1, e_2, \dots, e_K\}$.

③ The matrix product $(X E_{1:K}) E_{1:K}^t$ gives the best reconstruction of $X$ using only $K$ dimensions. $\left(\begin{array}{l} \text{ie: the projection of } X \text{ onto the} \\ \text{"best" } K\text{-dim subspace} \end{array}\right)$

④ The sum of eigenvalues $\sum_{j=1}^{K} \lambda_j$ gives the <u>total</u> variance of the projection of $X$ onto $E_{1:K}$

⑤ The ratio of sums

$$\frac{\sum_{j=1}^{K} \lambda_j}{\sum_{j=1}^{p} \lambda_j}$$ gives the <u>percentage of variance</u> preserved when projecting $X$ onto $E_{1:K}$.