

PCA -- In practice

- The basic goal of PCA is to reduce the dimensionality of the data. Thus, one usually chooses:

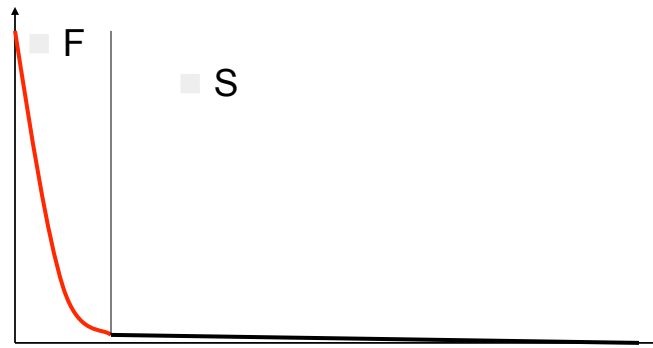
$$n \ll q$$

- But how do we select the number of components n ?

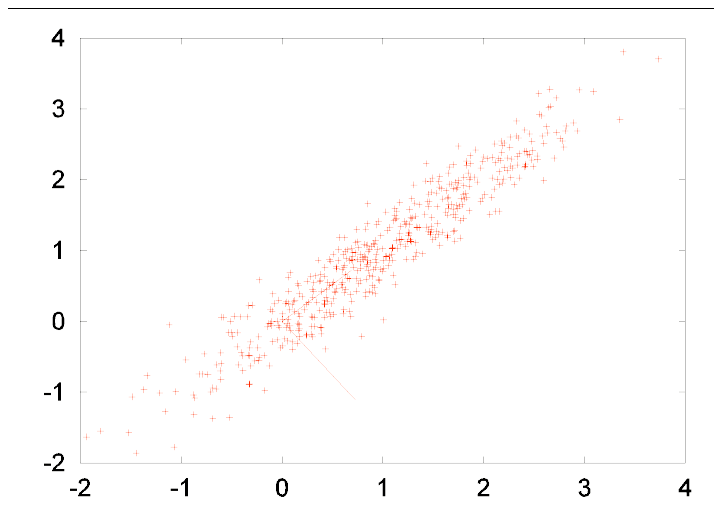
Determining the number of components

- Plot the eigenvalues – each eigenvalue is related to the amount of variation explained by the corresponding axis (eigenvector);
- If the points on the graph tend to level out (show an “elbow” shape), these eigenvalues are usually close enough to zero that they can be ignored.
- In general: Limit the variance accounted for.

Critical information lies in low dimensional subspaces



- A typical eigenvalue spectrum and its division into two orthogonal subspaces



$$\lambda_1 = 1.98, \lambda_2 = 0.05$$

Determining the number of components

$$\mathbf{x}_i \in \mathfrak{R}^q, \quad i = 1, \dots, N$$

$\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_q$: q eigenvectors (principal component directions)

$$\|\mathbf{w}_i\| = 1 \quad (\text{the } \mathbf{w}_i\text{s are orthonormal vectors})$$

Representation of \mathbf{x}_i in eigenvector space :

$$\mathbf{y}_i = (\mathbf{w}_1^T \mathbf{x}_i) \mathbf{w}_1 + (\mathbf{w}_2^T \mathbf{x}_i) \mathbf{w}_2 + \dots + (\mathbf{w}_q^T \mathbf{x}_i) \mathbf{w}_q$$

Suppose we retain the first k principal components :

$$\mathbf{y}_i^k = (\mathbf{w}_1^T \mathbf{x}_i) \mathbf{w}_1 + (\mathbf{w}_2^T \mathbf{x}_i) \mathbf{w}_2 + \dots + (\mathbf{w}_k^T \mathbf{x}_i) \mathbf{w}_k$$

Then :

$$\mathbf{y}_i - \mathbf{y}_i^k = (\mathbf{w}_{k+1}^T \mathbf{x}_i) \mathbf{w}_{k+1} + \dots + (\mathbf{w}_q^T \mathbf{x}_i) \mathbf{w}_q$$

Determining the number of components

$$(\mathbf{y}_i - \mathbf{y}_i^k)^T (\mathbf{y}_i - \mathbf{y}_i^k) =$$

$$\left[(\mathbf{w}_{k+1}^T \mathbf{x}_i) \mathbf{w}_{k+1} + \dots + (\mathbf{w}_q^T \mathbf{x}_i) \mathbf{w}_q \right]^T \left[(\mathbf{w}_{k+1}^T \mathbf{x}_i) \mathbf{w}_{k+1} + \dots + (\mathbf{w}_q^T \mathbf{x}_i) \mathbf{w}_q \right] =$$

$$\mathbf{w}_{k+1}^T (\mathbf{w}_{k+1}^T \mathbf{x}_i)^2 \mathbf{w}_{k+1} + \dots + \mathbf{w}_q^T (\mathbf{w}_q^T \mathbf{x}_i)^2 \mathbf{w}_q =$$

(note $\mathbf{w}_i^T \mathbf{w}_j = 0 \quad \forall i \neq j$ since \mathbf{w}_i and \mathbf{w}_j are orthogonal vectors)

$$(\mathbf{w}_{k+1}^T \mathbf{x}_i)^2 \mathbf{w}_{k+1}^T \mathbf{w}_{k+1} + \dots + (\mathbf{w}_q^T \mathbf{x}_i)^2 \mathbf{w}_q^T \mathbf{w}_q =$$

$$(\mathbf{w}_{k+1}^T \mathbf{x}_i)^2 + \dots + (\mathbf{w}_q^T \mathbf{x}_i)^2 =$$

$$(\mathbf{w}_{k+1}^T \mathbf{x}_i)(\mathbf{x}_i^T \mathbf{w}_{k+1}) + \dots + (\mathbf{w}_q^T \mathbf{x}_i)(\mathbf{x}_i^T \mathbf{w}_q) =$$

$$\mathbf{w}_{k+1}^T (\mathbf{x}_i \mathbf{x}_i^T) \mathbf{w}_{k+1} + \dots + \mathbf{w}_q^T (\mathbf{x}_i \mathbf{x}_i^T) \mathbf{w}_q$$

Determining the number of components

$$\frac{1}{N} \sum_{i=1}^N (\mathbf{y}_i - \mathbf{y}_i^k)^T (\mathbf{y}_i - \mathbf{y}_i^k) = \text{Mean square error}$$

$$\frac{1}{N} \sum_{i=1}^N [\mathbf{w}_{k+1}^T (\mathbf{x}_i \mathbf{x}_i^T) \mathbf{w}_{k+1} + \dots + \mathbf{w}_q^T (\mathbf{x}_i \mathbf{x}_i^T) \mathbf{w}_q] =$$

$$\mathbf{w}_{k+1}^T \left[\frac{1}{N} \sum_{i=1}^N (\mathbf{x}_i \mathbf{x}_i^T) \right] \mathbf{w}_{k+1} + \dots + \mathbf{w}_q^T \left[\frac{1}{N} \sum_{i=1}^N (\mathbf{x}_i \mathbf{x}_i^T) \right] \mathbf{w}_q =$$

$$\mathbf{w}_{k+1}^T \Sigma \mathbf{w}_{k+1} + \dots + \mathbf{w}_q^T \Sigma \mathbf{w}_q$$

$$\text{We have: } \Sigma \mathbf{w}_{k+1} = \lambda_{k+1} \mathbf{w}_{k+1}, \dots, \Sigma \mathbf{w}_q = \lambda_q \mathbf{w}_q$$

Thus:

$$\mathbf{w}_{k+1}^T \Sigma \mathbf{w}_{k+1} + \dots + \mathbf{w}_q^T \Sigma \mathbf{w}_q =$$

$$\mathbf{w}_{k+1}^T \lambda_{k+1} \mathbf{w}_{k+1} + \dots + \mathbf{w}_q^T \lambda_q \mathbf{w}_q =$$

$$\lambda_{k+1} + \dots + \lambda_q$$

Determining the number of components

$$\frac{1}{N} \sum_{i=1}^N (\mathbf{y}_i - \mathbf{y}_i^k)^T (\mathbf{y}_i - \mathbf{y}_i^k) = \lambda_{k+1} + \dots + \lambda_q$$



The mean square error of the truncated representation is equal to the sum of the remaining eigenvalues.

In general: choose k so that 90-95% of the variance of the data is captured.