

in the least squares sense, subject to the conditions $a_i^\top a_j = \delta_{ij}$, for all i, j with $1 \leq i, j \leq k$, where the matrix of the system is a block diagonal matrix consisting of k diagonal blocks $(X, \mathbf{1})$, where $\mathbf{1}$ denotes the column vector $(1, \dots, 1) \in \mathbb{R}^n$.

Again it is easy to see that each hyperplane H_i must pass through the centroid μ of X_1, \dots, X_n , and by switching to the centered data $X_i - \mu$ we get the system

$$\begin{pmatrix} X - \mu & 0 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & X - \mu \end{pmatrix} \begin{pmatrix} a_1 \\ \vdots \\ a_k \end{pmatrix} = \begin{pmatrix} 0 \\ \vdots \\ 0 \end{pmatrix},$$

with $a_i^\top a_j = \delta_{ij}$ for all i, j with $1 \leq i, j \leq k$.

If $VDU^\top = X - \mu$ is an SVD decomposition, it is easy to see that a least squares solution of this system is given by the last k columns of U , assuming that the main diagonal of D consists of the singular values $\sigma_1 \geq \sigma_2 \geq \cdots \geq \sigma_d$ of $X - \mu$ arranged in descending order. But now the $(d - k)$ -dimensional subspace U_{d-k} cut out by the hyperplanes defined by a_1, \dots, a_k is simply the orthogonal complement of (a_1, \dots, a_k) , which is the subspace spanned by the first $d - k$ columns of U .

So the best $(d - k)$ -dimensional affine subspace A_k approximating X_1, \dots, X_n in the least squares sense is

$$A_k = \mu + U_{d-k},$$

where U_{d-k} is the linear subspace spanned by the first $d - k$ principal directions of $X - \mu$, that is, the first $d - k$ columns of U . Consequently, we get the following interesting interpretation of PCA (actually, principal directions):

Theorem 23.12. *Let X be an $n \times d$ matrix of data points X_1, \dots, X_n , and let μ be the centroid of the X_i 's. If $X - \mu = VDU^\top$ is an SVD decomposition of $X - \mu$ and if the main diagonal of D consists of the singular values $\sigma_1 \geq \sigma_2 \geq \cdots \geq \sigma_d$, then a best $(d - k)$ -dimensional affine approximation A_k of X_1, \dots, X_n in the least squares sense is given by*

$$A_k = \mu + U_{d-k},$$

where U_{d-k} is the linear subspace spanned by the first $d - k$ columns of U , the first $d - k$ principal directions of $X - \mu$ ($1 \leq k \leq d - 1$).

Example 23.11. Going back to Example 23.10, a best 1-dimensional affine approximation A_1 is the affine line passing through $(\mu_1, \mu_2) = (1824.4, 5.6)$ of direction $u_1 = (0.9995, 0.0325)$.

Example 23.12. Suppose in the data set of Example 23.5 that we add the month of birth of every mathematician as a feature. We obtain the following data set.