**Chapter 2**

# Singular value decompositions

11 September 2016

©David Pollard 2016

# 1    SVD

Suppose once again that $X$ is an $n \times p$ matrix with $n \geq p$. It can be shown (Axler, 2015, Section 7.D) that there exist orthonormal bases $u_1, \dots, u_n$ for $\mathbb{R}^n$ and $v_1, \dots, v_p$ for $\mathbb{R}^p$, and numbers $\lambda_1 \geq \lambda_2, \cdots \geq \lambda_p \geq 0$, such that

$$Xv_j = \lambda_j u_j \qquad \text{for } 1 \leq j \leq p. \qquad\qquad <2.1>$$

In matrix form,

$$X(v_1, \dots, v_p) = (u_1, \dots, u_p)\mathrm{diag}(\lambda_1, \dots, \lambda_p) \qquad \text{or } XV = U\Lambda$$
$$\text{where } V = (v_1, \dots, v_p), \quad U = (u_1, \dots, u_p), \quad \Lambda = \mathrm{diag}(\lambda_1, \dots, \lambda_p).$$

Equivalently, because $V^T V = I_p = VV^T$,

$$X = U\Lambda V' = \sum\nolimits_{j \leq p} \lambda_j u_j v_j',$$

which is called the ***singular value decomposition*** of $X$. The $\lambda_i$'s are called the singular values of $X$.

If $\lambda_j > 0$ for $j \le m$ and zero otherwise then

$$u_j = \lambda_j^{-1} X v_j \qquad \text{for } 1 \le j \le m,$$

which shows that $u_1, \ldots, u_m$ all belong the the subspace $\mathcal{X}$ spanned by the columns of $X$. In the other direction, if $b = \sum_{j \le p} \theta_j v_j \in \mathbb{R}^p$ then

$$X b = \sum_{j \le p} \lambda_j u_j v_j' b = \sum_{j \le m} \lambda_j \langle v_j, b \rangle u_j = \sum_{j \le m} \lambda_j \theta_j u_j, \qquad \text{<2.2>}$$

which shows that every linear combination of the columns of $X$ can be written as a linear combination of $u_1, \ldots, u_m$. In other words, the subspace $\mathcal{X}$ has dimension $m$ and $\{u_1, \ldots, u_m\}$ is an orthonormal basis for $\mathcal{X}$.

The matrix $X$ has full rank if and only if $\lambda_j > 0$ for all $j$.

> **Remark.** We also have $X^T X = U \Lambda V^T V \Lambda U^T = U \Lambda^2 U^T$, which shows that the symmetric $p \times p$ matrix has eigenvalues $\lambda_1^2, \ldots, \lambda_p^2$ with corresponding eigenvectors $u_j$. That is, $X^T X u_j = \lambda_j^2 u_j$ for $1 \le j \le p$.

# 2     Least squares

If $\text{rank}(X) = m$, the matrix

$$H = (u_1, \ldots, u_m)(u_1, \ldots, u_m)^T = \sum_{j \le m} u_j u_j^T$$

projects vectors in $\mathbb{R}^n$ orthogonally onto $\mathcal{X}$. In particular, the orthogonal projection $\widehat{y}$ of a vector $y$ onto $\mathcal{X}$ is given by

$$\widehat{y} = H y = \sum_{i \le m} s_i u_i \qquad \text{where } s_i = \langle u_i, y \rangle.$$

If $X$ is of full rank then $\widehat{y}$ has a unique representation as $X\widehat{b}$, for some $\widehat{b}$ in $\mathbb{R}^p$. When $X$ has rank $m$, the SVD gives a neat expression for all possible solutions.

Write $y$ as $\sum_{i \le n} s_i u_i$, where $s_i = \langle u_i, y \rangle$, and write $b$ as $\sum_{j \le p} \theta_j v_j$, where $\theta_j = \langle v_j, b \rangle$. Then we need to find all values for $\theta$ such that

$$\sum_{i \le m} s_i u_i = \widehat{y} = X b = \sum_{j \le m} \lambda_j \theta_j u_j.$$

The solution is given by $\theta_j = s_j/\lambda_j$ for $1 \leq j \leq m$ and all other $\theta_j$'s uncon-strained. That is, $Xb = \widehat{y}$ if and only if

$$b = \sum_{j \leq m} (s_j/\lambda_j)v_j + w, \qquad \qquad <2.3>$$

where $w$ is an element of the $(p - m)$-dimensional subspace of $\mathbb{R}_p$ spanned by the unit vectors $v_{m+1}, \ldots, v_p$.

If $m$ equals $p$ (the full rank case) then $\widehat{b} = \sum_{j \leq p}(s_j/\lambda_j)v_j$ is the unique solution.

> **Remark.** If $X$ has rank $m$ then $Xv_j = \lambda_j u_j = 0$ for $j > m$. The $w$ in equation $<2.3>$ always contributes 0 to $Xb$.

# 3 Spectral norm (can be skipped)

The singular vector, $v_1$, is also the solution to a maximization problem: find the unit vector $t = \sum_j t_j v_j$ in $\mathbb{R}^p$ that maximizes $\|Xt\|$. Indeed

$$\|Xt\|^2 = \left\|\sum_{j \leq p} \lambda_j t_j u_j\right\|^2 = \sum_{j \leq m} \lambda_j^2 t_j^2,$$

a convex combination of $\lambda_1^2, \ldots, \lambda_p^2$ because $1 = \|t\|^2 = t_1^2 + \cdots + t_p^2$. The norm $\|Xt\|$ achieves its maximum value of $\lambda_1$ when $t_1^2 = 1$. The value $\lambda_1$ is often called the **spectral norm** of $X$, and is denoted by $\|X\|$ or $\|X\|_2$. It appears in many theoretical calculations and approximations involving $X$.

# 4 Perturbations and solutions of least squares problems (skip a little bit)

Suppose $X$ has full rank, so that the unique $\widehat{b}$ is given by $<2.3>$ with $w = 0$,

$$\widehat{b} = \sum_{j \leq p}(s_j/\lambda_j)v_j.$$

The representation suggests that if $\lambda_p$ is very small then it could have a large effect on $\widehat{b}$. This idea is not quite correct. After all, we could make $\lambda_p$ as large as we like by multiplying $X$ by a large enough constant, without really changing the least squares problem.

In fact, for $X$ of full rank, the relevant quantity is the ratio $\kappa(X) = \lambda_1/\lambda_p$, the so-called **condition number** of the matrix. When this ratio is large, relatively small numerical errors can be magnified into relatively large problems.

Dongarra et al. (1993, pages 9.5 and 11.4) gave a general bound for the effect of small perturbations of $X$ on the $\widehat{b}$. I'll settle for a much simpler construction that captures the main idea.

As before write a vector $y$ in $\mathbb{R}^n$ as $\sum_{i \leq n} s_i u_i$, with $s_i = \langle u_i, y \rangle$.

Define $E = u_p v_p^T$. You should convince yourself that $\|E\|_2 = 1$. For a small $\epsilon > 0$ define

$$X_\epsilon = X - \epsilon E = \left( \sum_{i < p} \lambda_i u_i v_i^T \right) + (\lambda_p - \epsilon) u_p v_p^T.$$

That is, the only change in the svd is reduction of the smallest singular value by $\epsilon$ (assuming $\epsilon < \lambda_p$). The relative size of the perturbation equals

$$\|X_\epsilon - X\|_2 / \|X\|_2 = \epsilon/\lambda_1.$$

The $\widehat{b}_\epsilon$ that minimizes $\|y - X_\epsilon b\|^2$ is given by

$$\widehat{b}_\epsilon = \left( \sum_{j < p} (s_j/\lambda_j) v_j \right) + v_p s_p/(\lambda_p - \epsilon) = \widehat{b} - s_p v_p \left( \lambda_p^{-1} - (\lambda_p - \epsilon)^{-1} \right).$$

That is,

$$\widehat{b}_\epsilon - \widehat{b} = (s_p/\lambda_p) \left( 1 - (1 - \epsilon/\lambda_p)^{-1} \right) v_p \approx \epsilon s_p/\lambda_p^2$$

so that $\left\| \widehat{b}_\epsilon - \widehat{b} \right\| \approx \epsilon |s_p|/\lambda_p^2$. Compare with

$$\left\| \widehat{b} \right\| = \sqrt{\sum_j s_j^2/\lambda_j^2} \approx |s_p|/\lambda_p$$

if $\lambda_p$ is much smaller than the other singular values and $|s_p|$ is not much smaller than the other $|s_j|$'s. In that case, the relative change in $\widehat{b}$ is approximately

$$\left\| \widehat{b}_\epsilon - \widehat{b} \right\| / \left\| \widehat{b} \right\| \approx \epsilon/\lambda_p.$$

and

$$\frac{\left\|\widehat{b}_\epsilon - \widehat{b}\right\| / \left\|\widehat{b}\right\|}{\left\|X_\epsilon - X\right\|_2 / \left\|X\right\|_2} \approx \lambda_1/\lambda_p = \kappa(X).$$

A lot of hand-waving in there, but it does suggest the magnifying effect of the condition number.

See the handout Longley.pdf for some calculations involving relatively small perturbations of a famous data set, which was used by Longley (1967) to discuss the effect of round-off errors on least squares. His abstract:

*Although there are many linear least squares programs available for use on the electronic computer, the algorithms specified in many of these programs are numerically more appropriate for the desk calculator than for the electronic computer. Routines which may be efficient for desk calculators may not be efficient for electronic computers. Since most computers carry about eight digits in the calculations, routines which do not take the problem of round-off errors and truncation into account may produce inaccurate numerical results. The difficulty is that the user will not know whether the results are accurate. Experiments with routine test problems using economic data indicated that either the data must be modified to fit the program or that the program must be altered to fit the data before numerical accuracy could be obtained on most programs tested. If the full potential of the electronic computer is to be achieved, an understanding of the basic arithmetic operations and their effect on the accuracy of the results is essential.*

In the handout I mess with $\widehat{b}$ by making small changes in $X$ and $y$.

# References

Axler, S. J. (2015). *Linear Algebra Done Right* (Third ed.). Undergraduate Texts in Mathematics. Springer.

Dongarra, J. J., C. B. Moler, J. R. Bunch, and G. W. Stewart (1993). *Linpack Users' Guide*. Society for Industrial and Applied Mathematics.

Longley, J. W. (1967). An appraisal of least squares programs for the electronic computer from the point of view of the user. *Journal of the American Statistical association 62*(319), 819–841.