

The singular value decomposition

While we have now derived the solution to the least squares problem (at least in certain special cases), our solution so far has not provided much insight into what kind of properties we can expect the solution to have. A particularly useful lens for thinking about least squares is through the **singular value decomposition** of the matrix \mathbf{A} in the objective function $\|\mathbf{y} - \mathbf{A}\mathbf{x}\|_2^2$. The singular value decomposition, or SVD for short, is just one of a family of useful matrix decompositions that you might encounter in data science. In general, a matrix decomposition is where we take a matrix \mathbf{A} and re-express it as a product (or sometimes a sum) of other simpler matrices that more clearly reveal some important structure of \mathbf{A} .

The SVD of a real-valued¹ $M \times N$ matrix \mathbf{A} is simply a factorization of the form

$$\mathbf{A} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T,$$

where \mathbf{U} , $\mathbf{\Sigma}$, and \mathbf{V} satisfy a number of properties, described below. Note that these properties involve a number of concepts from linear algebra such as the notion of an **orthnormal basis**, the **rank** of a matrix, and the concept of **eigenvectors** and **eigenvalues**. If these are a little fuzzy, review the linear algebra primer at the end of these notes before proceeding.

1. \mathbf{U} is an $M \times R$ matrix

$$\mathbf{U} = \begin{bmatrix} | & | & \cdots & | \\ \mathbf{u}_1 & \mathbf{u}_2 & \cdots & \mathbf{u}_R \\ | & | & & | \end{bmatrix}.$$

¹You can just as easily define the SVD for complex-valued matrices, but then every matrix transpose has to be replaced with a *Hermetian transpose*, in which you take both a transpose and compute the complex conjugate of all the entries. In this class we will only work with real-valued matrices, so we will stick with the simpler notation.

whose columns $\mathbf{u}_m \in \mathbb{R}^M$ are **orthonormal**. Note that while $\mathbf{U}^T \mathbf{U} = \mathbf{I}$, in general $\mathbf{U} \mathbf{U}^T \neq \mathbf{I}$ when $R < M$. The columns of \mathbf{U} are an orthobasis for the column space of \mathbf{A} . The columns of \mathbf{U} are called the **left-singular vectors** of \mathbf{A} .

2. \mathbf{V} is an $N \times R$ matrix

$$\mathbf{V} = \begin{bmatrix} | & | & \cdots & | \\ \mathbf{v}_1 & \mathbf{v}_2 & \cdots & \mathbf{v}_R \\ | & | & \cdots & | \end{bmatrix}.$$

whose columns $\mathbf{v}_n \in \mathbb{R}^N$ are orthonormal. Again, while $\mathbf{V}^T \mathbf{V} = \mathbf{I}$, in general $\mathbf{V} \mathbf{V}^T \neq \mathbf{I}$ when $R < N$. The columns of \mathbf{V} are an orthobasis for the row space of \mathbf{A} . The columns of \mathbf{V} are called the **right-singular vectors** of \mathbf{A} .

3. $\mathbf{\Sigma}$ is an $R \times R$ diagonal matrix with positive entries:

$$\mathbf{\Sigma} = \begin{bmatrix} \sigma_1 & 0 & \cdots & 0 \\ 0 & \sigma_2 & & \vdots \\ \vdots & & \ddots & \vdots \\ 0 & \cdots & \cdots & \sigma_R \end{bmatrix}.$$

We call the σ_r the **singular values** of \mathbf{A} . By convention, we will order them such that $\sigma_1 \geq \sigma_2 \geq \cdots \geq \sigma_R$.

What is important for our purposes is that *any* matrix (no matter its dimensions or entries) can be factorized in this form.

We can say a little bit more about how to interpret the SVD by recalling the notion of an **eigenvector** and its corresponding **eigenvalue**. Recall that for a square matrix \mathbf{B} , we call a vector \mathbf{x} an eigenvector if it satisfies

$$\mathbf{B}\mathbf{x} = \lambda\mathbf{x},$$

for some $\lambda \in \mathbb{C}$. We call λ the eigenvalue associated with \mathbf{x} .

Note that it is possible for \mathbf{B} to have complex eigenvalues, even if \mathbf{B} contains only real numbers. However, an important fact is that if \mathbf{B} is symmetric, i.e., if $\mathbf{B} = \mathbf{B}^T$, then all of its eigenvalues are *real*. An important class of matrices that we will arise frequently throughout this course are: **positive semidefinite matrices**. When we say that a matrix is positive semidefinite, we mean that it is symmetric and all of its eigenvalues are non-negative. If a symmetric matrix has eigenvalues that are strictly greater than zero, we call it **positive definite**.

There is a lot more that we could say about eigenvectors and eigenvalues, but for now all we want to point out is that the \mathbf{u}_m and the \mathbf{v}_n in the SVD can be interpreted as eigenvectors of matrices related to \mathbf{A} . For $M \neq N$, \mathbf{A} does not even have eigenvectors, but the \mathbf{u}_m are eigenvectors of $\mathbf{A}\mathbf{A}^T$, and the \mathbf{v}_n are eigenvectors of $\mathbf{A}^T\mathbf{A}$.

To see this, note that

$$\mathbf{A}\mathbf{A}^T = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T\mathbf{V}\mathbf{\Sigma}\mathbf{U}^T = \mathbf{U}\mathbf{\Sigma}^2\mathbf{U}^T,$$

since $\mathbf{V}^T\mathbf{V} = \mathbf{I}$. Note that since $\mathbf{\Sigma}$ is a diagonal matrix, $\mathbf{\Sigma}^2$ is just the same diagonal matrix, but where the entries along the diagonal are squared: σ_r^2 . Now consider $\mathbf{A}\mathbf{A}^T\mathbf{u}_m$. If we let \mathbf{e}_m denote the m^{th} “standard basis element”, i.e., the vector of all zeros with a single 1 in the m^{th} entry, then note that:

$$\begin{aligned}\mathbf{A}\mathbf{A}^T\mathbf{u}_m &= \mathbf{U}\mathbf{\Sigma}^2\mathbf{U}^T\mathbf{u}_m \\ &= \mathbf{U}\mathbf{\Sigma}^2\mathbf{e}_m \\ &= \mathbf{U}\sigma_m^2\mathbf{e}_m \\ &= \sigma_m^2\mathbf{u}_m.\end{aligned}$$

This shows that \mathbf{u}_m is an eigenvector of $\mathbf{A}\mathbf{A}^T$. Moreover, σ_m^2 is the corresponding eigenvalue. Thus, the singular values $\sigma_1, \dots, \sigma_R$ are the square roots of the non-zero eigenvalues of $\mathbf{A}\mathbf{A}^T$. This also shows that $\mathbf{A}\mathbf{A}^T$, in addition to being symmetric (since $(\mathbf{A}\mathbf{A}^T)^T = \mathbf{A}\mathbf{A}^T$), is in fact positive semidefinite, since σ_m^2 is always positive.

By essentially the same argument, you can show that $\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_R$ are eigenvectors of $\mathbf{A}^T\mathbf{A}$, and that the singular values $\sigma_1, \dots, \sigma_R$ are *also* the square roots of the non-zero eigenvalues of the positive semidefinite matrix $\mathbf{A}^T\mathbf{A}$.

The rank R is the dimension of the space spanned by the columns of \mathbf{A} , this is the same as the dimension of the space spanned by the rows. Thus $R \leq \min(M, N)$. We say \mathbf{A} is **full rank** if $R = \min(M, N)$.

When \mathbf{A} is **overdetermined** ($M > N$), the decomposition looks like this

$$\begin{bmatrix} \mathbf{A} \end{bmatrix} = \begin{bmatrix} \mathbf{U} \end{bmatrix} \begin{bmatrix} \sigma_1 & & \\ & \ddots & \\ & & \sigma_R \end{bmatrix} \begin{bmatrix} \mathbf{V}^T \end{bmatrix}.$$

When \mathbf{A} is **underdetermined** ($M < N$), the SVD looks like this

$$\begin{bmatrix} \mathbf{A} \end{bmatrix} = \begin{bmatrix} \mathbf{U} \end{bmatrix} \begin{bmatrix} \sigma_1 & & \\ & \ddots & \\ & & \sigma_R \end{bmatrix} \begin{bmatrix} \mathbf{V}^T \end{bmatrix}.$$

When \mathbf{A} is **square** and full rank ($M = N = R$), the SVD looks like

$$\begin{bmatrix} \mathbf{A} \end{bmatrix} = \begin{bmatrix} \mathbf{U} \end{bmatrix} \begin{bmatrix} \sigma_1 & & \\ & \ddots & \\ & & \sigma_N \end{bmatrix} \begin{bmatrix} \mathbf{V}^T \end{bmatrix}.$$

Alternative forms of the SVD

In most popular software packages, if you ask it to compute the SVD of a matrix \mathbf{A} , it will return matrices \mathbf{U} , $\mathbf{\Sigma}$, and \mathbf{V} , but the dimensions will differ from what I have written above. In particular, the default output will be:

- an $M \times M$ matrix $\tilde{\mathbf{U}}$,
- an $M \times N$ matrix $\tilde{\mathbf{\Sigma}}$,
- and $N \times N$ matrix $\tilde{\mathbf{V}}$.

So how does this output relate to what we have given above? Quite simply, \mathbf{U} is the first R columns of $\tilde{\mathbf{U}}$, $\mathbf{\Sigma}$ is the first R columns and rows of $\tilde{\mathbf{\Sigma}}$, and \mathbf{V} is the first R columns of $\tilde{\mathbf{V}}$.

This raises the question of what is happening in the remaining columns of $\tilde{\mathbf{U}}$, $\tilde{\mathbf{\Sigma}}$, and $\tilde{\mathbf{V}}$. For $\tilde{\mathbf{\Sigma}}$, the additional rows and columns are all zeros. In the case where $R < M$, so that $\mathbf{U} \neq \tilde{\mathbf{U}}$, we can write

$$\tilde{\mathbf{U}} = [\mathbf{U} \mid \mathbf{U}_0],$$

where \mathbf{U}_0 is an $M \times (M - R)$ matrix whose columns are an orthonormal basis for the part of \mathbb{R}^M that is orthogonal to the columns of \mathbf{U} . If $R < N$, so that $\mathbf{V} \neq \tilde{\mathbf{V}}$, we can similarly write

$$\tilde{\mathbf{V}} = [\mathbf{V} \mid \mathbf{V}_0],$$

where \mathbf{V}_0 is an $N \times (N - R)$ matrix whose columns are an orthonormal basis for the part of \mathbb{R}^N orthogonal to the columns of \mathbf{V} .

Note that, since $\tilde{\mathbf{\Sigma}}$ is just a zero-padded version of $\mathbf{\Sigma}$, it is not hard to show that

$$\tilde{\mathbf{U}}\tilde{\mathbf{\Sigma}}\tilde{\mathbf{V}}^T = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T.$$

The zeros in $\tilde{\Sigma}$ simply eliminate the contribution from \mathbf{U}_0 and \mathbf{V}_0 , so we are equally entitled to write $\mathbf{A} = \tilde{\mathbf{U}}\tilde{\Sigma}\tilde{\mathbf{V}}^T$. Just as was the case with \mathbf{U} and \mathbf{V} , we can write

$$\tilde{\mathbf{U}}^T\tilde{\mathbf{U}} = \mathbf{I} \quad \tilde{\mathbf{V}}^T\tilde{\mathbf{V}} = \mathbf{I}.$$

However, in this case, since $\tilde{\mathbf{U}}$ and $\tilde{\mathbf{V}}$ are square matrices, we also have that

$$\tilde{\mathbf{U}}\tilde{\mathbf{U}}^T = \mathbf{I} \quad \tilde{\mathbf{V}}\tilde{\mathbf{V}}^T = \mathbf{I}.$$

This is a consequence of a general fact about square matrices that you will prove on the homework. An equivalent way to think about this fact is that $\tilde{\mathbf{U}}$ is an orthonormal basis for all of \mathbb{R}^M and $\tilde{\mathbf{V}}$ is an orthonormal basis for all of \mathbb{R}^N . Please remember, however, that in general this is not true for \mathbf{U} and \mathbf{V} .

The SVD and least squares

So what does the SVD have to do with least squares? Recall that we had previously argued that if the matrix $\mathbf{A}^T\mathbf{A}$ is invertible, then the optimization problem

$$\underset{\mathbf{x} \in \mathbb{R}^N}{\text{minimize}} \quad \|\mathbf{y} - \mathbf{A}\mathbf{x}\|_2^2 \tag{1}$$

has solution

$$\hat{\mathbf{x}} = (\mathbf{A}^T\mathbf{A})^{-1}\mathbf{A}^T\mathbf{y}. \tag{2}$$

Now what can we say about this formula using the fact that we can write $\mathbf{A} = \mathbf{U}\Sigma\mathbf{V}^T$? First, note that

$$\mathbf{A}^T\mathbf{A} = \mathbf{V}\Sigma\mathbf{U}^T\mathbf{U}\Sigma\mathbf{V}^T = \mathbf{V}\Sigma^2\mathbf{V}^T.$$

Next, we use the fact that

$$(\mathbf{V}\Sigma^2\mathbf{V}^T)^{-1} = \mathbf{V}\Sigma^{-2}\mathbf{V}^T.$$

This is easy to check:

$$\mathbf{V}\Sigma^2\mathbf{V}^T(\mathbf{V}\Sigma^{-2}\mathbf{V}^T) = \mathbf{V}\Sigma^2\Sigma^{-2}\mathbf{V}^T = \mathbf{V}\mathbf{V}^T = \mathbf{I}.$$

Recall that, in general, we cannot always conclude that $\mathbf{V}\mathbf{V}^T = \mathbf{I}$. However, here we assume that the $N \times N$ matrix $\mathbf{A}^T\mathbf{A}$ is invertible. This matrix can only be invertible when it is full-rank, meaning that $R = N$, in which case, we also have $\mathbf{V}\mathbf{V}^T = \mathbf{I}$.

Returning to the least squares problem, using the above, we have that

$$\begin{aligned} (\mathbf{A}^T\mathbf{A})^{-1}\mathbf{A}^T &= \mathbf{V}\Sigma^{-2}\mathbf{V}^T\mathbf{V}\Sigma\mathbf{U}^T \\ &= \mathbf{V}\Sigma^{-2}\Sigma\mathbf{U}^T \\ &= \mathbf{V}\Sigma^{-1}\mathbf{U}^T. \end{aligned}$$

Thus, we arrive at an alternative way to express the solution to (1), now in terms of the SVD of \mathbf{A} :

$$\hat{\mathbf{x}} = \mathbf{V}\Sigma^{-1}\mathbf{U}^T\mathbf{y}. \quad (3)$$

This way of writing the least squares solution has a few nice advantages. First, recall that our previous formula in (2) only applied when $\mathbf{A}^T\mathbf{A}$ was invertible. It is not too hard to show (although we will not do so here) that (3) provides the solution to (1) regardless of whether $\mathbf{A}^T\mathbf{A}$ is invertible or not.

Moreover, it also provides the solution to a related problem. Suppose that we have a system $\mathbf{y} = \mathbf{A}\mathbf{x}$ which is *underdetermined*, meaning

that $M < N$. In this case, the least squares problem in (1) has infinitely many solutions. An alternative approach that might make more sense in this context (which can also be thought of as a kind of “least squares”) is to, from all \mathbf{x} that satisfy $\mathbf{y} = \mathbf{A}\mathbf{x}$, choose the one that is “smallest” in the Euclidean norm sense. As an optimization problem we could write this as:

$$\underset{\mathbf{x} \in \mathbb{R}^N}{\text{minimize}} \quad \|\mathbf{x}\|_2^2 \quad \text{subject to} \quad \mathbf{y} = \mathbf{A}\mathbf{x}. \quad (4)$$

It turns out that $\mathbf{V}\Sigma^{-1}\mathbf{U}^T\mathbf{y}$ is also the solution to (4).

Finally, it is relatively straightforward to show that in the case where \mathbf{A} is square and invertible ($M = N = R$), then $\mathbf{V}\Sigma^{-1}\mathbf{U}^T$ is simply \mathbf{A}^{-1} .

Least squares in noise

The real advantage to thinking about least squares through the lens of the SVD is that it tells us exactly what to expect when our measurements are corrupted with *noise*. Specifically, suppose we observe

$$\mathbf{y} = \mathbf{A}\mathbf{x} + \mathbf{e},$$

where $\mathbf{e} \in \mathbb{R}^M$ is an (unknown) perturbation of our measurements. Suppose that we then form the least squares estimate:

$$\hat{\mathbf{x}}_{\text{ls}} = \mathbf{A}^\dagger \mathbf{y} = \mathbf{A}^\dagger \mathbf{A}\mathbf{x} + \mathbf{A}^\dagger \mathbf{e}. \quad (5)$$

We would like to understand the effect of the noise vector \mathbf{e} on our estimate of \mathbf{x} . One way to understand this is to compare $\hat{\mathbf{x}}_{\text{ls}}$ to what

we would have obtained if we had been able to use least squares on the noise-free observations $\mathbf{y}_{\text{clean}} = \mathbf{A}\mathbf{x}$. The noise-free reconstruction is

$$\hat{\mathbf{x}}_{\text{clean}} = \mathbf{A}^\dagger \mathbf{y} = \mathbf{A}^\dagger \mathbf{A}\mathbf{x}.$$

Combining this with (5) we see that the additional error due to noise can be measured as

$$\|\hat{\mathbf{x}}_{\text{ls}} - \hat{\mathbf{x}}_{\text{clean}}\|_2^2 = \|\mathbf{A}^\dagger \mathbf{e}\|_2^2 = \|\mathbf{V}\boldsymbol{\Sigma}^{-1}\mathbf{U}^T \mathbf{e}\|_2^2. \quad (6)$$

Now suppose for a moment that the error has unit norm, $\|\mathbf{e}\|_2 = 1$. Just how bad can the error be? The worst case for (6) in this case is given by

$$\underset{\mathbf{e} \in \mathbb{R}^M}{\text{maximize}} \|\mathbf{V}\boldsymbol{\Sigma}^{-1}\mathbf{U}^T \mathbf{e}\|_2^2 \quad \text{subject to} \quad \|\mathbf{e}\|_2 = 1. \quad (7)$$

Note that $\boldsymbol{\Sigma}^{-1}\mathbf{U}^T \mathbf{e}$ is a vector in \mathbb{R}^R , and that for any vector $\mathbf{z} \in \mathbb{R}^R$, we have

$$\|\mathbf{V}\mathbf{z}\|_2^2 = \mathbf{z}^T \mathbf{V}^T \mathbf{V} \mathbf{z} = \mathbf{z}^T \mathbf{z} = \|\mathbf{z}\|_2^2,$$

since the columns of \mathbf{V} are orthonormal. Thus, we can simplify (7) to

$$\underset{\mathbf{e} \in \mathbb{R}^M}{\text{maximize}} \|\boldsymbol{\Sigma}^{-1}\mathbf{U}^T \mathbf{e}\|_2^2 \quad \text{subject to} \quad \|\mathbf{e}\|_2 = 1. \quad (8)$$

We can simplify a bit further by noticing that, since the columns of \mathbf{U} are orthonormal, $\|\mathbf{U}^T \mathbf{e}\|_2^2 \leq \|\mathbf{e}\|_2^2$. To see why this is the case, recall from our discussion of alternative forms of the SVD that when $R < M$, we can form the matrix

$$\tilde{\mathbf{U}} = [\mathbf{U} \mid \mathbf{U}_0],$$

where $\tilde{\mathbf{U}}$ is an $M \times M$ orthonormal matrix. It follows immediately that

$$\|\tilde{\mathbf{U}}^T \mathbf{e}\|_2^2 = \mathbf{e}^T \tilde{\mathbf{U}} \tilde{\mathbf{U}}^T \mathbf{e} = \mathbf{e}^T \mathbf{e} = \|\mathbf{e}\|_2^2,$$

since $\tilde{\mathbf{U}} \tilde{\mathbf{U}}^T = \mathbf{I}$. But it should also be clear that since $\tilde{\mathbf{U}}^T \mathbf{e}$ is just the vector of inner products between the columns in $\tilde{\mathbf{U}}$ and \mathbf{e} , and $\tilde{\mathbf{U}}$ is just the concatenation of \mathbf{U} and \mathbf{U}_0 ,

$$\|\tilde{\mathbf{U}}^T \mathbf{e}\|_2^2 = \|\mathbf{U}^T \mathbf{e}\|_2^2 + \|\mathbf{U}_0^T \mathbf{e}\|_2^2,$$

and hence we must have $\|\mathbf{U}^T \mathbf{e}\|_2^2 \leq \|\mathbf{e}\|_2^2$.

With this fact in hand, we can consider a slight “relaxation” of (8) to

$$\underset{\boldsymbol{\beta} \in \mathbb{R}^R}{\text{maximize}} \quad \|\boldsymbol{\Sigma}^{-1} \boldsymbol{\beta}\|_2^2 \quad \text{subject to} \quad \|\boldsymbol{\beta}\|_2 = 1.$$

We are not explicitly enforcing the fact that $\boldsymbol{\beta}$ should be a linear combination of the columns of \mathbf{U} here, which is why I am calling this a relaxation of the original problem. However, this problem has a simple solution that you will verify in the homework. Specifically the worst case $\boldsymbol{\beta}$ will have a 1 in the entry corresponding to the largest entry in $\boldsymbol{\Sigma}^{-1}$, and will be zero everywhere else. Thus

$$\max_{\boldsymbol{\beta} \in \mathbb{R}^R: \|\boldsymbol{\beta}\|_2=1} \|\boldsymbol{\Sigma}^{-1} \boldsymbol{\beta}\|_2^2 = \max_{r=1, \dots, R} \sigma_r^{-2} = \frac{1}{\sigma_R^2}.$$

(Recall that by convention, we order the singular values so that $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_R$.)

Note that, even though we ignored the constraint that $\boldsymbol{\beta}$ should be a linear combination of the columns of \mathbf{U} above, it is easy to find an \mathbf{e} such that $\boldsymbol{\beta} = \mathbf{U}^T \mathbf{e}$ gives us a 1 in the entry corresponding to σ_R .

In particular, $\mathbf{e} = \mathbf{u}_R$ gives us precisely this $\boldsymbol{\beta}$, and so even though we ignored this constraint, the result would not have changed had we included this requirement.

Returning to the reconstruction error (6), we now see that

$$\|\hat{\mathbf{x}}_{\text{ls}} - \hat{\mathbf{x}}_{\text{clean}}\|_2^2 = \|\mathbf{V}\boldsymbol{\Sigma}^{-1}\mathbf{U}^T\mathbf{e}\|_2^2 \leq \frac{1}{\sigma_R^2}\|\mathbf{e}\|_2^2.$$

While this is a worst-case bound, notice that if σ_R is small, the worst case reconstruction error can be **very bad**. Later in this course we will discuss some strategies to mitigate this fact.

Linear Algebra Review II: Inner products, orthonormal bases, and eigenvectors

We have already encountered the notion of a general vector space with an abstract norm. One other key concept, which we have already treated informally, is the notion of an abstract *inner product*.

Inner products

Definition: An **inner product** on a real-valued² vector space \mathcal{S} is a mapping

$$\langle \cdot, \cdot \rangle : \mathcal{S} \times \mathcal{S} \rightarrow \mathbb{R}$$

that obeys

1. $\langle \mathbf{x}, \mathbf{y} \rangle = \langle \mathbf{y}, \mathbf{x} \rangle$

2. For any $a, b \in \mathbb{R}$

$$\langle a\mathbf{x} + b\mathbf{y}, \mathbf{z} \rangle = a\langle \mathbf{x}, \mathbf{z} \rangle + b\langle \mathbf{y}, \mathbf{z} \rangle$$

3. $\langle \mathbf{x}, \mathbf{x} \rangle \geq 0$ and $\langle \mathbf{x}, \mathbf{x} \rangle = 0 \Leftrightarrow \mathbf{x} = \mathbf{0}$

Standard Examples:

- $\mathcal{S} = \mathbb{R}^N$,

$$\langle \mathbf{x}, \mathbf{y} \rangle = \sum_{n=1}^N x_n y_n = \mathbf{y}^T \mathbf{x}$$

²By real-valued, we mean the scalar associated with scalar multiplication is a real-number. You can easily extend these definitions to a complex-valued vector space, but since we will not need this level of generality in this course we will stick with the simpler real-valued case.

- $\mathcal{S} = L_2([a, b])$,

$$\langle \mathbf{x}, \mathbf{y} \rangle = \int_a^b x(t)y(t) dt$$

Slightly less standard example:

- $\mathcal{S} = \mathbb{R}^{M \times N}$ (the set of $M \times N$ matrices with real entries)

$$\langle \mathbf{X}, \mathbf{Y} \rangle = \text{trace}(\mathbf{Y}^T \mathbf{X}) = \sum_{m=1}^M \sum_{n=1}^N X_{m,n} Y_{m,n}$$

(Recall that $\text{trace}(\mathbf{X})$ is the sum of the entries on the diagonal of \mathbf{X} .) This is called the *trace inner product* or *Frobenius inner product* or *Hilbert-Schmidt inner product*.

A vector space equipped with an inner product is called an **inner product space**. Inner products allow us to think about angles between vectors in arbitrary vector spaces – most importantly, we can generalize the notion of orthogonality from Euclidean space to an arbitrary inner product space: two vectors \mathbf{x}, \mathbf{y} are **orthogonal** if $\langle \mathbf{x}, \mathbf{y} \rangle = 0$.

Induced norms

A valid inner product induces a valid norm by

$$\|\mathbf{x}\| = \sqrt{\langle \mathbf{x}, \mathbf{x} \rangle}$$

(You can check that $\sqrt{\langle \mathbf{x}, \mathbf{x} \rangle}$ indeed satisfies the properties required of a norm on your own as an exercise.)

It is not hard to see that in \mathbb{R}^N , the standard inner product induces the ℓ_2 norm, i.e., $\sqrt{\mathbf{x}^T \mathbf{x}} = \|\mathbf{x}\|_2$.

Properties of induced norms

In addition to the triangle inequality,

$$\|\mathbf{x} + \mathbf{y}\| \leq \|\mathbf{x}\| + \|\mathbf{y}\|,$$

which all norms must obey, induced norms obey some very handy inequalities. Below I give two of the most famous and useful ones. Note that these are not necessarily true for norms in general, only for norms induced by an inner product:

- **Pythagorean Theorem**

$$\langle \mathbf{x}, \mathbf{y} \rangle = 0 \Rightarrow \|\mathbf{x} + \mathbf{y}\|^2 = \|\mathbf{x}\|^2 + \|\mathbf{y}\|^2$$

The left-hand side above also implies that $\|\mathbf{x} - \mathbf{y}\|^2 = \|\mathbf{x}\|^2 + \|\mathbf{y}\|^2$.

- **Cauchy-Schwarz Inequality**

$$|\langle \mathbf{x}, \mathbf{y} \rangle| \leq \|\mathbf{x}\| \|\mathbf{y}\|$$

Equality is achieved above when (and only when) \mathbf{x} and \mathbf{y} are **colinear**:

$$\exists a \in \mathbb{R} \quad \text{such that} \quad \mathbf{y} = a\mathbf{x}.$$

The last remaining set of concepts from linear algebra that we need to review relate to the notion of a **basis** (and if we have an inner product, an **orthonormal basis**). To get offer formal definitions of these terms, we first need to review a few basic ideas.

Linear combinations and spans

Let $\mathcal{M} = \{\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_N\}$ be a set of vectors in a real-valued vector space \mathcal{S} .

Definition: A **linear combination** of vectors in \mathcal{M} is a sum of the form

$$a_1\mathbf{v}_1 + a_2\mathbf{v}_2 + \cdots + a_N\mathbf{v}_N$$

for some $a_1, a_2, \dots, a_N \in \mathbb{R}$.

Definition: The **span** of \mathcal{M} is the set of all linear combinations of \mathcal{M} . We write this as

$$\text{span}(\mathcal{M}) = \text{span}(\{\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_N\})$$

Example:

Suppose

$$\mathcal{S} = \mathbb{R}^3, \quad \mathbf{v}_1 = \begin{bmatrix} 1 \\ 1 \\ 0 \end{bmatrix}, \quad \mathbf{v}_2 = \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix}.$$

In this case,

$$\text{span}(\{\mathbf{v}_1, \mathbf{v}_2\}) = \text{the } (x_1, x_2) \text{ plane.}$$

In other words, the span of $\mathbf{v}_1, \mathbf{v}_2$ is the set of any $\mathbf{x} \in \mathbb{R}^3$ of the form

$$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ 0 \end{bmatrix}.$$

It should be clear that for any such \mathbf{x} , we can write \mathbf{x} as a linear combination of \mathbf{v}_1 and \mathbf{v}_2 :

$$\mathbf{x} = x_1 \mathbf{v}_1 + (x_2 - x_1) \mathbf{v}_2.$$

Linear dependence

A set of vectors $\{\mathbf{v}_j\}_{j=1}^N$ is said to be **linearly dependent** if there exists scalars a_1, a_2, \dots, a_N , not all equal to 0, such that

$$\sum_{n=1}^N a_n \mathbf{v}_n = \mathbf{0}.$$

Note that in such a case, we can write (at least) one of the vectors \mathbf{v}_n as a linear combination of the others.

Likewise, if $\sum_n a_n \mathbf{v}_n = \mathbf{0}$ only when all the $a_j = 0$, then $\{\mathbf{v}_n\}_{n=1}^N$ is said to be **linearly independent**.

Example:

$$\mathcal{S} = \mathbb{R}^3, \quad \mathbf{v}_1 = \begin{bmatrix} 2 \\ 1 \\ 0 \end{bmatrix} \quad \mathbf{v}_2 = \begin{bmatrix} 1 \\ 1 \\ 0 \end{bmatrix} \quad \mathbf{v}_3 = \begin{bmatrix} 1 \\ 2 \\ 0 \end{bmatrix}$$

It is too hard to find an a_1, a_2, a_3 such that

$$a_1 \mathbf{v}_1 + a_2 \mathbf{v}_2 + a_3 \mathbf{v}_3 = \mathbf{0}.$$

Note that any two of the vectors above are linearly independent:

$$\begin{aligned} \text{span}(\{\mathbf{v}_1, \mathbf{v}_2, \mathbf{v}_3\}) &= \text{span}(\{\mathbf{v}_1, \mathbf{v}_2\}) \\ &= \text{span}(\{\mathbf{v}_1, \mathbf{v}_3\}) \\ &= \text{span}(\{\mathbf{v}_2, \mathbf{v}_3\}). \end{aligned}$$

Bases

Definition: A **basis** of a finite-dimensional³ vector space \mathcal{S} is a set of vectors \mathcal{B} such that

1. $\text{span}(\mathcal{B}) = \mathcal{S}$
2. \mathcal{B} is linearly independent

The second condition ensures that all bases of \mathcal{S} will have the same number of elements.

The **dimension** of \mathcal{S} is the number of elements required in a basis for \mathcal{S} .

Examples:

1. \mathbb{R}^N with

$$\{\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_N\} = \left\{ \begin{bmatrix} 1 \\ 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix}, \dots, \begin{bmatrix} 0 \\ 0 \\ 0 \\ \vdots \\ 1 \end{bmatrix} \right\}$$

This is the **standard basis** for \mathbb{R}^N .

The dimension of \mathbb{R}^N is N .

2. \mathbb{R}^N with any set of N linearly independent vectors.

3. $\mathcal{S} = \{\text{polynomials of degree at most } p\}$.

A basis for \mathcal{S} is $\mathcal{B} = \{1, t, t^2, \dots, t^p\}$.

The dimension of \mathcal{S} is $p + 1$.

³Things get a bit more complicated in infinite-dimensional vector spaces. We won't worry about those details in this class.

Column space, row space, and rank

Consider the matrix \mathbf{V} which has the vectors $\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_N \in \mathbb{R}^M$ as columns:

$$\mathbf{V} = \left[\begin{array}{c|c|c|c} | & | & \cdots & | \\ \mathbf{v}_1 & \mathbf{v}_2 & \cdots & \mathbf{v}_n \\ | & | & \cdots & | \end{array} \right].$$

Note that this is an $M \times N$ matrix.

The **column space**, also commonly called the **range**, of \mathbf{V} is just the span of the columns of \mathbf{V} , i.e., $\text{span}(\{\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_N\})$. This is often denoted by $\mathcal{R}(\mathbf{V})$. Note that $\mathcal{R}(\mathbf{V}) \subseteq \mathbb{R}^M$ (since the \mathbf{v}_j are vectors in \mathbb{R}^M). It may be possible to have $\mathcal{R}(\mathbf{V}) = \mathbb{R}^M$ – this will occur if \mathbf{V} has M linearly independent columns. Note that \mathbf{V} cannot have more than M linearly independent columns: since \mathbb{R}^M is an M -dimensional space, you cannot have more than M linearly independent vectors within that space. Also, note that if $N < M$, then the columns of \mathbf{V} cannot possibly span all of \mathbb{R}^M . Moreover, even if $N > M$, it is still possible that $\mathcal{R}(\mathbf{V})$ is a strict subspace of \mathbb{R}^M .

The **row space** is defined similarly, and is simply the span of the *rows* of \mathbf{V} , or equivalently, the row space is the column space of \mathbf{V}^T . Note that the row space is a subspace of \mathbb{R}^N .

It is a fact that for any matrix \mathbf{V} , the column space and row space have the same dimension – that is, the maximum number of linearly independent columns is the same as the maximum number of linearly independent rows. This dimension is called the **rank** of \mathbf{V} .

Orthogonal bases

A collection of vectors $\{\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_N\}$ in a finite dimensional vector space \mathcal{S} is called an **orthogonal basis** if

1. $\text{span}(\{\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_N\}) = \mathcal{S}$,
2. $\langle \mathbf{v}_j, \mathbf{v}_k \rangle = 0$ for all $j \neq k$.

If in addition the vectors are normalized (under the induced norm),

$$\|\mathbf{v}_n\| = 1, \quad \text{for } n = 1, \dots, N,$$

we will call it an **orthonormal basis** or **orthobasis**.

In the case where $\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_N$ are vectors in \mathbb{R}^M , we can form an $M \times N$ matrix \mathbf{V} with them as columns:

$$\mathbf{V} = \begin{bmatrix} | & | & \cdots & | \\ \mathbf{v}_1 & \mathbf{v}_2 & \cdots & \mathbf{v}_N \\ | & | & \cdots & | \end{bmatrix}.$$

In this case, another way to express that $\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_N$ are orthonormal is that

$$\mathbf{V}^T \mathbf{V} = \mathbf{I}.$$

To see this, note that $\mathbf{V}^T \mathbf{v}_j$ computes the inner product between \mathbf{v}_j and all N vectors in the basis:

$$\begin{bmatrix} - & \mathbf{v}_1^T & - \\ & \vdots & \\ - & \mathbf{v}_N^T & - \end{bmatrix} \begin{bmatrix} | \\ \mathbf{v}_j \\ | \end{bmatrix} = \begin{bmatrix} \langle \mathbf{v}_j, \mathbf{v}_1 \rangle \\ \vdots \\ \langle \mathbf{v}_j, \mathbf{v}_N \rangle \end{bmatrix}.$$

These inner products will all be zero except for the j^{th} element. Repeating this for $\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_N$ results in $\mathbf{V}^T \mathbf{V} = \mathbf{I}$.