

Stable Reconstruction with the Truncated SVD

We have seen that if \mathbf{A} has very small singular values and we apply the pseudo-inverse in the presence of noise, the results can be disastrous. But it doesn't have to be this way. There are several ways to stabilize the pseudo-inverse. We start by discussing the simplest one, where we simply “cut out” the part of the reconstruction which is causing the problems.

As before, we are given noisy indirect observations of a vector \mathbf{x} through a $M \times N$ matrix \mathbf{A} :

$$\mathbf{y} = \mathbf{A}\mathbf{x} + \mathbf{e}. \quad (1)$$

The matrix \mathbf{A} has SVD $\mathbf{A} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T$, and pseudo-inverse $\mathbf{A}^\dagger = \mathbf{V}\mathbf{\Sigma}^{-1}\mathbf{U}^T$.

At this point it is useful to recall that we could write the matrix $\mathbf{U}\mathbf{V}^T$ as a sum of outer products of the columns of \mathbf{U} and \mathbf{V} :

$$\mathbf{U}\mathbf{V}^T = \sum_{r=1}^R \mathbf{u}_r \mathbf{v}_r^T,$$

where R is the rank of \mathbf{A} , and $\mathbf{u}_r \in \mathbb{R}^M$ and $\mathbf{v}_r \in \mathbb{R}^N$ are columns of \mathbf{U} and \mathbf{V} , respectively. See the addendum at the end of these notes on matrix multiplication if this way of writing a matrix product seems unfamiliar.

Since $\mathbf{\Sigma}$ is diagonal, we can think of $\mathbf{U}\mathbf{\Sigma}$ as just rescaling the columns of \mathbf{U} , so that we can also rewrite $\mathbf{A} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T$ as a sum:

$$\mathbf{A} = \sum_{r=1}^R \sigma_r \mathbf{u}_r \mathbf{v}_r^T,$$

where the σ_r are the singular values. Similarly, we can write the pseudo-inverse as

$$\mathbf{A}^\dagger = \sum_{r=1}^R \frac{1}{\sigma_r} \mathbf{v}_r \mathbf{u}_r^\top.$$

Given \mathbf{y} as above, we can write the least-squares estimate of \mathbf{x} from the noisy measurements as

$$\hat{\mathbf{x}}_{\text{ls}} = \mathbf{A}^\dagger \mathbf{y} = \sum_{r=1}^R \frac{1}{\sigma_r} \langle \mathbf{y}, \mathbf{u}_r \rangle \mathbf{v}_r. \quad (2)$$

As we can see (and have seen before) if any one of the σ_r are very small, the least squares reconstruction can be a disaster.

A simple way to avoid this is to simply truncate the sum (2), leaving out the terms where σ_r is too small ($1/\sigma_r$ is too big). Exactly how many terms to keep depends a great deal on the application, as there are competing interests. On the one hand, we want to ensure that each of the σ_r we include has an inverse of reasonable size, on the other, we want the reconstruction to be accurate (i.e., not to deviate from the noiseless least squares solution by too much).

We form an approximation \mathbf{A}' to \mathbf{A} by taking

$$\mathbf{A}' = \sum_{r=1}^{R'} \sigma_r \mathbf{u}_r \mathbf{v}_r^\top,$$

for some $R' < R$. Again, our final answer will depend on which R' we use, and choosing R' is often times something of an art. It is clear that the approximation \mathbf{A}' has rank R' . Note that the pseudo-inverse of \mathbf{A}' is also a truncated sum

$$\mathbf{A}'^\dagger = \sum_{r=1}^{R'} \frac{1}{\sigma_r} \mathbf{v}_r \mathbf{u}_r^\top.$$

Given noisy data \mathbf{y} as in (1), we reconstruct \mathbf{x} by applying the truncated pseudo-inverse to \mathbf{y} :

$$\hat{\mathbf{x}}_{\text{trunc}} = \mathbf{A}'^\dagger \mathbf{y} = \sum_{r=1}^{R'} \frac{1}{\sigma_r} \langle \mathbf{y}, \mathbf{u}_r \rangle \mathbf{v}_r.$$

How good is this reconstruction? To answer this question, we will again compare it to the least squares reconstruction corresponding to “noiseless” measurements $\hat{\mathbf{x}}_{\text{clean}} = \mathbf{A}^\dagger \mathbf{A} \mathbf{x}$. The difference between these two is the reconstruction error (relative to $\hat{\mathbf{x}}_{\text{clean}}$) as

$$\begin{aligned} \hat{\mathbf{x}}_{\text{trunc}} - \hat{\mathbf{x}}_{\text{clean}} &= \mathbf{A}'^\dagger \mathbf{y} - \mathbf{A}^\dagger \mathbf{A} \mathbf{x} \\ &= \mathbf{A}'^\dagger \mathbf{A} \mathbf{x} + \mathbf{A}'^\dagger \mathbf{e} - \mathbf{A}^\dagger \mathbf{A} \mathbf{x} \\ &= (\mathbf{A}'^\dagger - \mathbf{A}^\dagger) \mathbf{A} \mathbf{x} + \mathbf{A}'^\dagger \mathbf{e}. \end{aligned}$$

Proceeding further, we can write the matrix $\mathbf{A}'^\dagger - \mathbf{A}^\dagger$ as

$$\mathbf{A}'^\dagger - \mathbf{A}^\dagger = \sum_{r=R'+1}^R -\frac{1}{\sigma_r} \mathbf{v}_r \mathbf{u}_r^\top,$$

and so the first term in the reconstruction error can be written as

$$\begin{aligned} (\mathbf{A}'^\dagger - \mathbf{A}^\dagger) \mathbf{A} \mathbf{x} &= \sum_{r=R'+1}^R -\frac{1}{\sigma_r} \langle \mathbf{A} \mathbf{x}, \mathbf{u}_r \rangle \mathbf{v}_r \\ &= \sum_{r=R'+1}^R -\frac{1}{\sigma_r} \left\langle \sum_{j=1}^R \sigma_j \langle \mathbf{x}, \mathbf{v}_j \rangle \mathbf{u}_j, \mathbf{u}_r \right\rangle \mathbf{v}_r \\ &= \sum_{r=R'+1}^R -\frac{1}{\sigma_r} \sum_{j=1}^R \sigma_j \langle \mathbf{x}, \mathbf{v}_j \rangle \langle \mathbf{u}_j, \mathbf{u}_r \rangle \mathbf{v}_r \\ &= \sum_{r=R'+1}^R -\langle \mathbf{x}, \mathbf{v}_r \rangle \mathbf{v}_r \quad (\text{since } \langle \mathbf{u}_r, \mathbf{u}_j \rangle = 0 \text{ unless } j = r). \end{aligned}$$

The second term in the reconstruction error can also be expanded against the \mathbf{v}_r :

$$\mathbf{A}'^\dagger \mathbf{e} = \sum_{r=1}^{R'} \frac{1}{\sigma_r} \langle \mathbf{e}, \mathbf{u}_r \rangle \mathbf{v}_r.$$

Combining these expressions, the reconstruction error can be written

$$\begin{aligned} \widehat{\mathbf{x}}_{\text{trunc}} - \widehat{\mathbf{x}}_{\text{clean}} &= \underbrace{\sum_{r=1}^{R'} \frac{1}{\sigma_r} \langle \mathbf{e}, \mathbf{u}_r \rangle \mathbf{v}_r}_{\text{Noise error}} + \underbrace{\sum_{r=R'+1}^R -\langle \mathbf{x}, \mathbf{v}_r \rangle \mathbf{v}_r}_{\text{Approximation error}} \\ &= \text{Noise error} + \text{Approximation error}. \end{aligned}$$

Since the \mathbf{v}_r are mutually orthogonal, and the two sums run over disjoint index sets, the noise error and the approximation error will be orthogonal. Also

$$\begin{aligned} \|\widehat{\mathbf{x}}_{\text{trunc}} - \widehat{\mathbf{x}}_{\text{clean}}\|_2^2 &= \|\text{Noise error}\|_2^2 + \|\text{Approximation error}\|_2^2 \\ &= \sum_{r=1}^{R'} \frac{1}{\sigma_r^2} |\langle \mathbf{e}, \mathbf{u}_r \rangle|^2 + \sum_{r=R'+1}^R |\langle \mathbf{x}, \mathbf{v}_r \rangle|^2. \end{aligned}$$

The reconstruction error, then, is signal dependent and will depend on how much of the vector \mathbf{x} is concentrated in the subspace spanned by $\mathbf{v}_{R'+1}, \dots, \mathbf{v}_R$. We will lose everything in this subspace. On the other hand, if it contains a significant part of \mathbf{x} , then there is not much least squares can do for you.

The worst-case noise error occurs when \mathbf{e} is aligned with $\mathbf{u}_{R'}$:

$$\|\text{Noise error}\|_2^2 = \sum_{r=1}^{R'} \frac{1}{\sigma_r^2} |\langle \mathbf{e}, \mathbf{u}_r \rangle|^2 \leq \frac{1}{\sigma_{R'}^2} \cdot \|\mathbf{e}\|_2^2.$$

Stable Reconstruction using Tikhonov Regularization

Tikhonov¹ regularization is another way to stabilize the least squares recovery. It has the nice features that: 1) it can be interpreted using optimization, and 2) it can be computed without direct knowledge of the SVD of \mathbf{A} .

Recall that we motivated the pseudo-inverse by showing that $\hat{\mathbf{x}}_{\text{ls}} = \mathbf{A}^\dagger \mathbf{y}$ is a solution to

$$\underset{\mathbf{x} \in \mathbb{R}^N}{\text{minimize}} \quad \|\mathbf{y} - \mathbf{A}\mathbf{x}\|_2^2. \quad (3)$$

When \mathbf{A} has full column rank, $\hat{\mathbf{x}}_{\text{ls}}$ is the unique solution, otherwise it is the solution with smallest energy. When \mathbf{A} has full column rank but has singular values which are very small, huge variations in \mathbf{x} (in directions of the singular vectors \mathbf{v}_r corresponding to the tiny σ_r) can have very little effect on the residual $\|\mathbf{y} - \mathbf{A}\mathbf{x}\|_2^2$. As such, the solution to (3) can have wildly inaccurate components in the presence of even mild noise.

One way to counteract this problem is to modify (3) with a **regularization** term that penalizes the size of the solution $\|\mathbf{x}\|_2^2$ as well as the residual error $\|\mathbf{y} - \mathbf{A}\mathbf{x}\|_2^2$:

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

The parameter $\delta > 0$ gives us a trade-off between accuracy and regularization; we want to choose δ small enough so that the residual for the solution of (4) is close to that of (3), and large enough so that the problem is well-conditioned (i.e., stable in the presence of noise).

¹Andrey Tikhonov (1906-1993) was a 20th century Russian mathematician.

Just as with (3), which is solved by applying the pseudo-inverse to \mathbf{y} , we can write the solution to (4) in closed form.

We have two ways of writing this solution. In terms of the SVD of \mathbf{A} , the minimizer of (4) is

$$\begin{aligned}\hat{\mathbf{x}}_{\text{tik}} &= \mathbf{V}\hat{\boldsymbol{\alpha}}_{\text{tik}} \\ &= \mathbf{V}(\boldsymbol{\Sigma}^2 + \delta\mathbf{I})^{-1}\boldsymbol{\Sigma}\mathbf{U}^T\mathbf{y}.\end{aligned}\tag{5}$$

Alternatively, we can write the solution in terms of \mathbf{A} as

$$\hat{\mathbf{x}}_{\text{tik}} = (\mathbf{A}^T\mathbf{A} + \delta\mathbf{I})^{-1}\mathbf{A}^T\mathbf{y}.\tag{6}$$

Note that even if $\mathbf{A}^T\mathbf{A}$ is not invertible $\mathbf{A}^T\mathbf{A} + \delta\mathbf{I}$ always is (if $\delta > 0$). Thus, the expression (6) holds for all M, N , and R . It is also the case that

$$\hat{\mathbf{x}}_{\text{tik}} = \mathbf{A}^T(\mathbf{A}\mathbf{A}^T + \delta\mathbf{I})^{-1}\mathbf{y}.$$

You will show that all of these expressions are equivalent on the next homework.

Tikhonov regularization is in some sense very similar to the truncated SVD, but with one significant advantage: because of the second way of writing the solution, we do not need to explicitly calculate the SVD to solve (4). The importance of not needing to explicitly compute the SVD is significant when we are solving large problems. When \mathbf{A} is large ($M, N > 10^5$, say) it may be expensive or even impossible to construct the SVD and compute with it explicitly. However, if it has special structure (if it is sparse, for example), then it may take many fewer than MN operations to compute a matrix vector product $\mathbf{A}\mathbf{x}$.

In these situations, a **matrix free** iterative algorithm can be used to perform the inverse required in (6). A prominent example of such

an algorithm is **gradient descent** and its many variants, which we will see very soon.

We can get a better feel for what Tikhonov regularization is doing by comparing it directly to the pseudo-inverse. Recall that the least squares reconstruction $\hat{\mathbf{x}}_{\text{ls}}$ can be written as

$$\hat{\mathbf{x}}_{\text{ls}} = \sum_{r=1}^R \frac{1}{\sigma_r} \langle \mathbf{y}, \mathbf{u}_r \rangle \mathbf{v}_r.$$

The Tikhonov reconstruction $\hat{\mathbf{x}}_{\text{tik}}$ derived above is

$$\hat{\mathbf{x}}_{\text{tik}} = \sum_{r=1}^R \frac{\sigma_r}{\sigma_r^2 + \delta} \langle \mathbf{y}, \mathbf{u}_r \rangle \mathbf{v}_r. \quad (7)$$

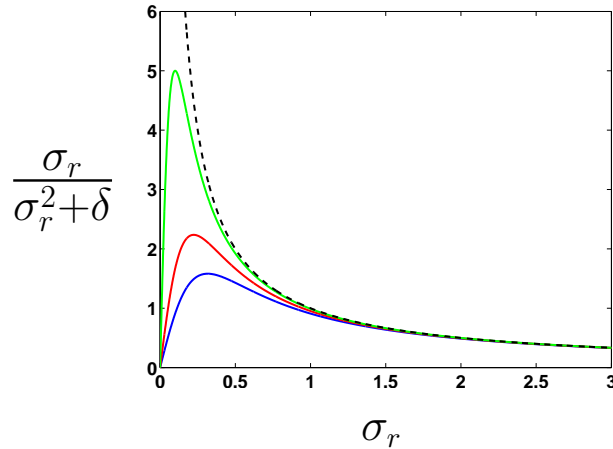
Notice that when σ_r is much larger than δ ,

$$\frac{\sigma_r}{\sigma_r^2 + \delta} \approx \frac{1}{\sigma_r}, \quad \sigma_r \gg \delta,$$

but when σ_r is small

$$\frac{\sigma_r}{\sigma_r^2 + \delta} \approx 0, \quad \sigma_r \ll \delta.$$

Thus the Tikhonov reconstruction modifies the important parts (components where the σ_r are large) of the pseudo-inverse very little, while ensuring that the unimportant parts (components where the σ_r are small) affect the solution only by a very small amount. This **damping** of the singular values, is illustrated below.



Above, we see the damped multipliers $\sigma_r/(\sigma_r^2 + \delta)$ versus σ_r for $\delta = 0.1$ (blue), $\delta = 0.05$ (red), and $\delta = 0.01$ (green). The black dotted line is $1/\sigma_r$, the least squares multiplier. Notice that for large σ_r ($\sigma_r > 2\sqrt{\delta}$, say), the damping has almost no effect.

This damping makes the Tikhonov reconstruction exceptionally stable; large multipliers never appear in the reconstruction (7). In fact it is easy to check that

$$\frac{\sigma_r}{\sigma_r^2 + \delta} \leq \frac{1}{2\sqrt{\delta}}$$

no matter the value of σ_r .

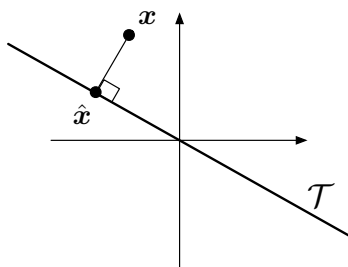
You can perform a very similar kind of noise analysis for Tikhonov reconstruction as we just did for the truncated SVD, but we will leave you to do this at home.

A geometric perspective on least squares

Before we finally move on to talk a bit more about practical algorithms, let's see one more way in which the SVD can help us to think about what least squares is doing.

Specifically, another way to think about least squares involves a simple geometric problem. Suppose we are given a vector $\mathbf{x} \in \mathbb{R}^M$ and a subspace² $\mathcal{T} = \text{span}(\{\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_N\})$ (where $N < M$).

If \mathbf{x} does not already live in the span of $\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_N$, we might ask “what is the closest point $\hat{\mathbf{x}} \in \mathcal{T}$ to \mathbf{x} ? This is illustrated below:



Mathematically, we want to find the $\hat{\mathbf{x}} \in \mathcal{T}$ that minimizes $\|\mathbf{x} - \hat{\mathbf{x}}\|_2$, i.e., given \mathbf{x} , we want to solve the optimization program

$$\underset{\mathbf{y} \in \mathbb{R}^M}{\text{minimize}} \quad \|\mathbf{x} - \mathbf{y}\|_2^2 \quad \text{subject to} \quad \mathbf{y} \in \mathcal{T}.$$

This might initially seem different than what we've been considering so far – it looks like a least squares problem, but we have the constraint that \mathbf{y} lives in a subspace. However, think for a moment about what the constraint $\mathbf{y} \in \mathcal{T}$ actually means.

²Recall that a subspace of \mathbb{R}^M is just a set of vectors that can be thought of as a vector space in its own right, meaning that any linear combination of vectors in the subspace produces another vector in that same subspace. The span of $N < M$ vectors in \mathbb{R}^M , such as a 2D plane in 3 dimensions, is the canonical example of a subspace.

If $\mathbf{y} \in \mathcal{T}$, then \mathbf{y} can be written as a linear combination of the vectors $\mathbf{a}_1, \dots, \mathbf{a}_N$. If \mathbf{A} is the $M \times N$ matrix with columns given by $\mathbf{a}_1, \dots, \mathbf{a}_N$, this is equivalent to saying that we can write $\mathbf{y} = \mathbf{A}\boldsymbol{\alpha}$ for some $\boldsymbol{\alpha} \in \mathbb{R}^N$. Thus, instead of optimizing over $\mathbf{y} \in \mathcal{T}$ we can simply optimize over $\boldsymbol{\alpha}$ in the equivalent problem

$$\underset{\boldsymbol{\alpha} \in \mathbb{R}^N}{\text{minimize}} \quad \|\mathbf{x} - \mathbf{A}\boldsymbol{\alpha}\|_2^2.$$

Note that this is simply the standard least squares problem we have been studying for a while now, with solution $\hat{\boldsymbol{\alpha}} = \mathbf{V}\boldsymbol{\Sigma}^{-1}\mathbf{U}^T\mathbf{x}$. This gives us a solution for $\hat{\mathbf{x}}$ of our original problem of

$$\begin{aligned} \hat{\mathbf{x}} &= \mathbf{A}\hat{\boldsymbol{\alpha}} \\ &= \mathbf{U}\boldsymbol{\Sigma}\mathbf{V}^T\mathbf{V}\boldsymbol{\Sigma}^{-1}\mathbf{U}^T\mathbf{x} \\ &= \mathbf{U}\boldsymbol{\Sigma}\boldsymbol{\Sigma}^{-1}\mathbf{U}^T\mathbf{x} \\ &= \mathbf{U}\mathbf{U}^T\mathbf{x}. \end{aligned}$$

Note that in the picture on the previous page, I indicated that the vector $\hat{\mathbf{x}} - \mathbf{x}$ was orthogonal to the line that represents the subspace \mathcal{T} . This is a general fact for least squares solutions: at the optimal $\hat{\mathbf{x}}$, we always have that the error $\hat{\mathbf{x}} - \mathbf{x}$ is orthogonal to $\mathcal{T} = \mathcal{R}(\mathbf{A})$.

This is easy to see using the SVD. Since for any $\mathbf{y} \in \mathcal{R}(\mathbf{A})$ we can write $\mathbf{y} = \mathbf{U}\boldsymbol{\Sigma}\mathbf{V}^T\boldsymbol{\alpha}$, we have

$$\begin{aligned} \langle \hat{\mathbf{x}} - \mathbf{x}, \mathbf{y} \rangle &= \langle \mathbf{U}\mathbf{U}^T\mathbf{x} - \mathbf{x}, \mathbf{U}\boldsymbol{\Sigma}\mathbf{V}^T\boldsymbol{\alpha} \rangle \\ &= \boldsymbol{\alpha}^T\mathbf{V}\boldsymbol{\Sigma}\mathbf{U}^T\mathbf{U}\mathbf{U}^T\mathbf{x} - \boldsymbol{\alpha}^T\mathbf{V}\boldsymbol{\Sigma}\mathbf{U}^T\mathbf{x} \\ &= \boldsymbol{\alpha}^T\mathbf{V}\boldsymbol{\Sigma}\mathbf{U}^T\mathbf{x} - \boldsymbol{\alpha}^T\mathbf{V}\boldsymbol{\Sigma}\mathbf{U}^T\mathbf{x} \\ &= 0, \end{aligned}$$

for any $\boldsymbol{\alpha} \in \mathbb{R}^N$.

The fact that the optimal approximation to \mathbf{x} in \mathcal{T} , $\hat{\mathbf{x}}$ results in an error $\hat{\mathbf{x}} - \mathbf{x}$ that is orthogonal to \mathcal{T} is often called the **orthogonality principle**. In fact, you can actually start a discussion of least squares by first showing that the orthogonality principle must be true, and then using this fact to derive the solution to the least squares problem.

Another piece of terminology you may encounter is that $\hat{\mathbf{x}}$ is the **orthogonal projection** of \mathbf{x} onto \mathcal{T} . In general, when people talk about an orthogonal projection, this is simply shorthand for taking a vector and finding the closest point in some subspace \mathcal{T} to \mathbf{x} . The fact that you do this by finding a point in \mathcal{T} from which the direction to \mathbf{x} is orthogonal to \mathcal{T} is why it gets the name “orthogonal”.

A note on matrix multiplication

Although we have encountered some of these concepts already, now is a good time to review some basic notions involving how we think about matrix multiplication. When dealing with simple matrix-vector products, these ideas are straightforward, but as we get to increasingly complex matrix factorizations, it can be harder to break down what is happening.

Matrix-vector multiplication

Let's begin by considering simple matrix-vector multiplication $\mathbf{A}\mathbf{x}$. There are really two ways to think about this. First, suppose that \mathbf{A} is an $M \times N$ matrix. We can think of \mathbf{A} as the concatenation of N columns, denoted by $\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_N$:

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

The first way of thinking of a matrix-vector multiplication is that $\mathbf{A}\mathbf{x}$ is a weighted combination of the columns of \mathbf{A} :

$$\begin{bmatrix} | & | & \cdots & | \\ \mathbf{a}_1 & \mathbf{a}_2 & \cdots & \mathbf{a}_N \\ | & | & \cdots & | \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_N \end{bmatrix} = x_1\mathbf{a}_1 + x_2\mathbf{a}_2 + \cdots + x_N\mathbf{a}_N.$$

To describe the second way of thinking of a matrix-vector multiplication, we are going to slightly abuse our notation. Even though we just said \mathbf{a}_n was a column of \mathbf{A} , which is what we will do for the rest

of the course, note that we can also break up a matrix by along its rows, which we will (just for now) denote by $\mathbf{a}_1^T, \mathbf{a}_2^T, \dots, \mathbf{a}_M^T$:

$$\mathbf{A} = \begin{bmatrix} \text{---} & \mathbf{a}_1^T & \text{---} \\ \text{---} & \mathbf{a}_2^T & \text{---} \\ & \vdots & \\ \text{---} & \mathbf{a}_M^T & \text{---} \end{bmatrix}.$$

We use a transpose when referring to the rows since normally \mathbf{a}_m always refers to a column vector. This way of writing \mathbf{A} suggests an alternative way to think about matrix-vector multiplication. That is, each entry of $\mathbf{A}\mathbf{x}$ is the inner product between the rows of \mathbf{A} and the vector \mathbf{x} :

$$\begin{bmatrix} \text{---} & \mathbf{a}_1^T & \text{---} \\ \text{---} & \mathbf{a}_2^T & \text{---} \\ & \vdots & \\ \text{---} & \mathbf{a}_M^T & \text{---} \end{bmatrix} \begin{bmatrix} | \\ | \\ | \\ | \\ | \\ \mathbf{x} \\ | \\ | \\ | \end{bmatrix} = \begin{bmatrix} \mathbf{a}_1^T \mathbf{x} \\ \mathbf{a}_2^T \mathbf{x} \\ \vdots \\ \mathbf{a}_M^T \mathbf{x} \end{bmatrix}.$$

Matrix-matrix multiplication

Likewise, the product of an $M \times N$ matrix \mathbf{A} and a $N \times P$ matrix \mathbf{B} can be thought of in two different ways. The traditional way (how you would compute the elements of \mathbf{AB} one entry at a time) is as a collection of the inner products between all of the rows of \mathbf{A} and all of the columns of \mathbf{B} . If we continue to be a little loose with our notation, we can visualize this as:

$$\begin{bmatrix} \text{---} & \mathbf{a}_1^T & \text{---} \\ \text{---} & \mathbf{a}_2^T & \text{---} \\ & \vdots & \\ \text{---} & \mathbf{a}_M^T & \text{---} \end{bmatrix} \begin{bmatrix} | & | & \cdots & | \\ \mathbf{b}_1 & \mathbf{b}_2 & \cdots & \mathbf{b}_P \\ | & | & & | \end{bmatrix} = \begin{bmatrix} \mathbf{a}_1^T \mathbf{b}_1 & \mathbf{a}_1^T \mathbf{b}_2 & \cdots & \mathbf{a}_1^T \mathbf{b}_P \\ \mathbf{a}_2^T \mathbf{b}_1 & \mathbf{a}_2^T \mathbf{b}_2 & \cdots & \mathbf{a}_2^T \mathbf{b}_P \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{a}_M^T \mathbf{b}_1 & \mathbf{a}_M^T \mathbf{b}_2 & \cdots & \mathbf{a}_M^T \mathbf{b}_P \end{bmatrix}.$$

However, just as with matrix-vector multiplication, there are alternative perspectives. We think of each column of \mathbf{AB} as being the result of a matrix-vector product that combines the columns of \mathbf{A} :

$$\mathbf{A} \left[\begin{array}{c|c|c|c} \mathbf{b}_1 & \mathbf{b}_2 & \cdots & \mathbf{b}_P \end{array} \right] = \left[\begin{array}{c|c|c|c} \mathbf{Ab}_1 & \mathbf{Ab}_2 & \cdots & \mathbf{Ab}_P \end{array} \right].$$

Similarly, we can also think of each row of \mathbf{AB} as a linear combination of the rows of \mathbf{B} :

$$\left[\begin{array}{c|c|c} \text{---} & \mathbf{a}_1^T & \text{---} \\ \text{---} & \mathbf{a}_2^T & \text{---} \\ & \vdots & \\ \text{---} & \mathbf{a}_M^T & \text{---} \end{array} \right] \mathbf{B} = \left[\begin{array}{c|c|c} \text{---} & \mathbf{a}_1^T \mathbf{B} & \text{---} \\ \text{---} & \mathbf{a}_2^T \mathbf{B} & \text{---} \\ & \vdots & \\ \text{---} & \mathbf{a}_M^T \mathbf{B} & \text{---} \end{array} \right].$$

Finally, there is a fourth way to think of matrix-matrix multiplication, as a sum of the rank 1 matrices formed by taking the outer product of the columns of \mathbf{A} with the rows of \mathbf{B} :

$$\left[\begin{array}{c|c|c|c} \mathbf{a}_1 & \mathbf{a}_2 & \cdots & \mathbf{a}_N \end{array} \right] \left[\begin{array}{c|c|c} \text{---} & \mathbf{b}_1^T & \text{---} \\ \text{---} & \mathbf{b}_2^T & \text{---} \\ & \vdots & \\ \text{---} & \mathbf{b}_N^T & \text{---} \end{array} \right] = \sum_{n=1}^N \mathbf{a}_n \mathbf{b}_n^T.$$

To see this, consider just \mathbf{a}_1 and what it contributes to \mathbf{AB} . Each column of \mathbf{AB} will potentially have some contribution from \mathbf{a}_1 . Just how much? Well, it's the first row of \mathbf{B} that will determine this for each column. The same argument follows for each of the N columns of \mathbf{A} and rows of \mathbf{B} .