

**CS189 Introduction to Machine Learning**

Spring 2018

Notes: Benjamin Recht

**Some notes on pseudoinverses**

Let  $\mathbf{A}$  be a square, full rank,  $d \times d$  dimensional matrix. We all know that the inverse of  $\mathbf{A}$  is the matrix  $\mathbf{B}$  such that  $\mathbf{AB} = \mathbf{I}$ . What's the natural notion of an "inverse" for a matrix that is not square or possibly not full rank?

Let  $\mathbf{X}$  be a  $n \times d$  dimensional matrix and assume that  $n \geq d$ . We say that a  $d \times n$  dimensional  $\mathbf{Z}$  is a *pseudoinverse* of  $\mathbf{X}$  if the following properties hold

1.  $\mathbf{XZ}\mathbf{X} = \mathbf{X}$ .
2.  $\mathbf{ZXZ} = \mathbf{Z}$ .
3.  $\mathbf{XZ}$  and  $\mathbf{ZX}$  are both symmetric matrices.

We denote the pseudoinverse by  $\mathbf{X}^\dagger$ .

Note that all of these properties are trivially satisfied by the matrix inverse. So for a square matrix with full rank  $\mathbf{A}^\dagger = \mathbf{A}^{-1}$ .

Let's first assume that  $\mathbf{X}$  is full rank. Then the  $\mathbf{X}^\dagger = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T$ . To verify this, observe

$$\mathbf{X} \mathbf{X}^\dagger \mathbf{X} = \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{X} = \mathbf{X}.$$

Similarly,

$$\mathbf{X}^\dagger \mathbf{X} \mathbf{X}^\dagger = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T = \mathbf{X}^\dagger.$$

Finally, we can see that

$$\begin{aligned} \mathbf{X}^\dagger \mathbf{X} &= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{X} = \mathbf{I} \\ \mathbf{X} \mathbf{X}^\dagger &= \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T, \end{aligned}$$

which are both symmetric matrices. Interestingly, in the case that  $\mathbf{X}$  is full rank, multiplying  $\mathbf{X}$  by the pseudoinverse on the left yields the identity matrix, just as a normal inverse would. On the other hand, since  $\mathbf{X} \mathbf{X}^\dagger$  is  $n \times n$  it is necessarily rank deficient. In this case,  $\mathbf{X} \mathbf{X}^\dagger$  is the projection onto the range of  $\mathbf{X}$ .

What about when the matrix is rank deficient? The pseudoinverse still exists and we can derive it using the singular value decomposition. Let

$$\mathbf{X} = \mathbf{U} \mathbf{S} \mathbf{V}^T = \sum_{i=1}^d \sigma_i \mathbf{u}_i \mathbf{v}_i^T$$

be a singular value decomposition of  $\mathbf{X}$  and

$$\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_d$$

be the singular values. Let  $k$  denote the index of the *smallest* nonzero singular value. That is, for  $j > k$ , we should have  $\sigma_j = 0$ . If  $\mathbf{X}$  is full rank, then  $k = d$ .

Now define

$$\mathbf{Z} := \sum_{j=1}^k \sigma_j^{-1} \mathbf{v}_j \mathbf{u}_j^T.$$

I claim that this matrix  $\mathbf{Z}$  satisfies all of the properties of a pseudoinverse. To see this, let's just work with matrices. Define the  $d \times n$  matrix  $\mathbf{R}$  to be

$$\mathbf{R} = \begin{bmatrix} \text{diag}(\sigma_1^{-1}, \sigma_2^{-1}, \dots, \sigma_k^{-1}, 0, \dots, 0) & \mathbf{0}_{d \times (n-d)} \end{bmatrix}.$$

Note that  $\mathbf{R}$  is the pseudoinverse of  $\mathbf{S}$ ! Moreover,

$$\begin{aligned} \mathbf{XZ}\mathbf{X} &= \mathbf{USV}^T \mathbf{VRU}^T \mathbf{USV}^T \\ &= \mathbf{USRSV}^T = \mathbf{USV}^T. \end{aligned}$$

The other identities can be similarly checked. Now, in this case, the one particularly interesting identity is

$$\mathbf{ZX} = \mathbf{VRSV}^T = \sum_{i=1}^k \mathbf{v}_i \mathbf{v}_i^T$$

In this case,  $\mathbf{X}^\dagger \mathbf{X}$  does not equal the identity matrix, but is instead the projection operator onto the range of  $\mathbf{X}^T$ .

Finally, consider the context of least-squares. In the case that  $\mathbf{X}$  is rank deficient, the least-squares solution is not unique. However, we can define a unique solution in the range of  $\mathbf{X}^T$  using the pseudoinverse. If we define  $\mathbf{w}_{\text{pinv}} = \mathbf{X}^\dagger \mathbf{y}$  this corresponds to the vector

$$\mathbf{w}_{\text{pinv}} = \sum_{i=1}^k \sigma_i^{-1} \langle \mathbf{u}_i, \mathbf{y} \rangle \mathbf{v}_i$$

where  $k$  again denotes the index of the smallest, nonzero singular value. This is precisely the least-squares solution when we restrict the possible  $\mathbf{w}$  to lie in the span of the data (i.e., the range of  $\mathbf{X}^T$ ).