

# Thoughts about covariance

Andy Smith

September 2012

We have  $N$  measurements stored in an array  $\mathbf{X}$ :

$$\mathbf{X} = [ \mathbf{x}_1 \ \mathbf{x}_2 \ \dots \ \mathbf{x}_N ] \quad (1)$$

where  $\mathbf{x}_i$  is the  $i$ -th measurement vector, and has  $M$  elements. If we say that the mean vector is

$$\bar{\mathbf{x}} = \frac{1}{N} \sum_{i=1}^N \mathbf{x}_i, \quad (2)$$

and will also have  $M$  elements, then we can define

$$\mathbf{Y} = [ \bar{\mathbf{x}} \ \bar{\mathbf{x}} \ \dots \ \bar{\mathbf{x}} ] = \bar{\mathbf{x}} \mathbf{1}_N^T \quad (3)$$

where  $\mathbf{1}_N$  is a column vector of length  $N$ , with every element set to 1.

Using  $\mathbf{X}$  and  $\mathbf{Y}$  we can write the covariance,  $\mathbf{S}$  as:

$$\mathbf{S} = \frac{1}{N-1} (\mathbf{X} - \mathbf{Y}) (\mathbf{X} - \mathbf{Y})^T \quad (4)$$

$$= \frac{1}{N-1} (\mathbf{X} - \bar{\mathbf{x}} \mathbf{1}_N^T) (\mathbf{X} - \bar{\mathbf{x}} \mathbf{1}_N^T)^T \quad (5)$$

$$= \frac{1}{N-1} \left( \mathbf{X}\mathbf{X}^T - (\mathbf{X}\mathbf{1}_N \bar{\mathbf{x}}^T)^T - \mathbf{X}\mathbf{1}_N \bar{\mathbf{x}}^T + \bar{\mathbf{x}} \mathbf{1}_N^T \mathbf{1}_N \bar{\mathbf{x}}^T \right). \quad (6)$$

But  $\mathbf{X}\mathbf{1}_N = N\bar{\mathbf{x}}$  (since it is the sum of elements in the rows of  $\mathbf{X}$ ), and  $\mathbf{1}_N^T \mathbf{1}_N = N$ , so that:

$$\mathbf{S} = \frac{1}{N-1} (\mathbf{X}\mathbf{X}^T - 2N\bar{\mathbf{x}}\bar{\mathbf{x}}^T + N\bar{\mathbf{x}}\bar{\mathbf{x}}^T) \quad (7)$$

$$\mathbf{S} = \frac{1}{N-1} (\mathbf{X}\mathbf{X}^T - N\bar{\mathbf{x}}\bar{\mathbf{x}}^T). \quad (8)$$

Now imagine that we want to combine two areas of a dataset together:

$$\mathbf{X}_1 = [ \mathbf{x}_1 \ \mathbf{x}_2 \ \dots \ \mathbf{x}_{N_1} ] \quad \text{and} \quad \mathbf{X}_2 = [ \mathbf{x}_{N_1+1} \ \mathbf{x}_{N_1+2} \ \dots \ \mathbf{x}_{N_2} ] \quad (9)$$

so that  $\mathbf{X} = [ \mathbf{X}_1 \ \mathbf{X}_2 ]$  gives us the original dataset. In this case, the mean vector can be obtained easily from:

$$\bar{\mathbf{x}} = \frac{N_1 \bar{\mathbf{x}}_1 + N_2 \bar{\mathbf{x}}_2}{N_1 + N_2}. \quad (10)$$

The vast majority of computing time for  $\mathbf{S}$  is in the matrix multiplication  $\mathbf{X}\mathbf{X}^T$ , so a way to combine  $\mathbf{X}_1\mathbf{X}_1^T$  and  $\mathbf{X}_2\mathbf{X}_2^T$  in a timely fashion would be very pleasing. Happily, since  $(\mathbf{X}\mathbf{X}^T)_{ij} = \sum_k x_{ik} x_{jk}$ , then  $\mathbf{X}\mathbf{X}^T = \mathbf{X}_1\mathbf{X}_1^T + \mathbf{X}_2\mathbf{X}_2^T$ .