

# Theano Computations

Jonathan Simon

January 25, 2016

## MLP Activation Layer

The hidden Layer in Theano's MLP is computed as:  $X \rightarrow f(XW + b) = A$ , implying that the activation vectors are rows, and the batch-wide unit activation are columns. For an activation layer with  $d$  hidden units and a minibatch of size  $t$ , the activation layer is given by the matrix:

$$A = \begin{bmatrix} a_{11} & \dots & a_{1d} \\ \vdots & & \vdots \\ a_{t1} & \dots & a_{td} \end{bmatrix} = \begin{bmatrix} | & & | \\ u_1 & \dots & u_d \\ | & & | \end{bmatrix} = \begin{bmatrix} \text{---} & o_1 & \text{---} \\ & \vdots & \\ \text{---} & o_t & \text{---} \end{bmatrix} \quad (1)$$

where:

$u_i = t$ -dimensional timecourse of the  $i^{\text{th}}$  unit

$o_i = d$ -dimensional state of the activation layer at the  $i^{\text{th}}$  observation

## Covariance Computation

We can then compute the between-unit covariance using the formulas provided on the wikipedia pages:

1. [https://en.wikipedia.org/wiki/Sample\\_mean\\_and\\_covariance](https://en.wikipedia.org/wiki/Sample_mean_and_covariance)
2. [https://en.wikipedia.org/wiki/Estimation\\_of\\_covariance\\_matrices](https://en.wikipedia.org/wiki/Estimation_of_covariance_matrices)
3. <https://en.wikipedia.org/wiki/Covariance>

Consider the definition of covariance for a pair of random variables  $U_j$  and  $U_k$  representing the idealized scalar activations of the  $j^{\text{th}}$  and  $k^{\text{th}}$  hidden units respectively:

$$\text{cov}(U_j, U_k) = \text{E}[(U_j - \text{E}[U_j])(U_k - \text{E}[U_k])] \quad (2)$$

Substituting the **random variables**  $U_j$  and  $U_k$  for the **column vectors**  $u_j$  and  $u_k$  of scalar activation observations:

$$\text{cov}(u_j, u_k) = \frac{1}{t-1} \left( (u_j - \mathbb{1}\bar{u}_j)^\top (u_k - \mathbb{1}\bar{u}_k) \right) \quad (3)$$

where the **scalar**  $\bar{u}_*$  is defined as:

$$\bar{u}_* = \frac{1}{t} \sum_{i=1}^t [u_*]_i \quad \left( = \frac{1}{t} \sum_{i=1}^t [o_i]_* \right)$$

We can then extend this calculation to all pairs of units in the hidden layer to compute the covariance matrix:

$$\Sigma = \frac{1}{t-1} \left( (A - \mathbb{1}\bar{o})^\top (A - \mathbb{1}\bar{o}) \right) \quad (4)$$

where the **row vector**  $\bar{o}$  is defined as:

$$\bar{o} = \frac{1}{t} \sum_{i=1}^t o_i$$

## Correlation Computation

We can also compute the between-unit pearson correlation using the formulas provided at:

1. <http://www.johndcook.com/blog/2008/11/05/how-to-calculate-pearson-correlation-accurately/>
2. [https://en.wikipedia.org/wiki/Pearson\\_product-moment\\_correlation\\_coefficient](https://en.wikipedia.org/wiki/Pearson_product-moment_correlation_coefficient)

Consider the definition of correlation for a pair of random variables  $U_j$  and  $U_k$  representing the idealized scalar activations of the  $j^{\text{th}}$  and  $k^{\text{th}}$  hidden units respectively:

$$\rho_{jk} = \frac{\text{cov}(U_j, U_k)}{\sigma_j \sigma_k} = \frac{\text{E}[(U_j - \text{E}[U_j])(U_k - \text{E}[U_k])]}{\sqrt{\text{E}[(U_j - \text{E}[U_j])^2]} \sqrt{\text{E}[(U_k - \text{E}[U_k])^2]}} \quad (5)$$

Substituting the **random variables**  $U_j$  and  $U_k$  for the **column vectors**  $u_j$  and  $u_k$  of scalar activation observations:

$$r_{jk} = \frac{(u_j - \mathbb{1}\bar{u}_j)^\top (u_k - \mathbb{1}\bar{u}_k)}{\sqrt{(u_j - \mathbb{1}\bar{u}_j)^\top (u_j - \mathbb{1}\bar{u}_j)} \sqrt{(u_k - \mathbb{1}\bar{u}_k)^\top (u_k - \mathbb{1}\bar{u}_k)}} \quad (6)$$

where the **scalar**  $\bar{u}_*$  is defined as:

$$\bar{u}_* = \frac{1}{t} \sum_{i=1}^t [u_*]_i \quad \left( = \frac{1}{t} \sum_{i=1}^t [o_i]_* \right)$$

Note that the normalizing factor of  $\frac{1}{t-1}$  is absent in the above expression because it is canceled by the two occurrences of  $\sqrt{\frac{1}{t-1}}$  in the denominator.

Define the sample variance  $s_*$  as:

$$s_* = \sqrt{\frac{1}{t-1} (u_* - \mathbb{1}\bar{u}_*)^\top (u_* - \mathbb{1}\bar{u}_*)} \quad (7)$$

Then we can compute the **row vector**  $s$  as:

$$s = [s_1 \quad \dots \quad s_d] = \sqrt{\frac{1}{t-1} \left( \mathbb{1}^\top \left( (A - \mathbb{1}\bar{o}) \odot (A - \mathbb{1}\bar{o}) \right) \right)} \quad (8)$$

Using the definition of the covariance matrix from the previous section, we can then write the correlation matrix as:

$$\rho = \frac{\Sigma}{(\mathbb{1}s) \odot (\mathbb{1}s)^\top} = \frac{\Sigma}{(\mathbb{1}s)^\top \odot (\mathbb{1}s)} \quad (9)$$

where the division is performed element-wise.

Additional optimizations which can be performed in the Theano code include:

1. Removing the  $\frac{1}{t-1}$  scaling factors, since they are present in both numerator and denominator
2. Removing the square roots, since the final step involves element-wise squaring
3. Replacing the multiplications by  $\mathbb{1}$  with broadcasting and axis-wise summation
4. Replacing the Hadamard products with element-wise squaring

**Update:** Implemented optimization 3 and 4, but optimizations 1 and 2 don't speed things up enough to make it worth the confusion.