

# Assignment 2

Wendy Nieuwkamer

## 1 Question 1

This question is about *vectorization*, i.e. writing expressions in matrix-vector form. The goal is to vectorize the update rule for multivariate linear regression.

Let  $\theta$  be the parameter vector  $\theta = (\theta_0 \quad \theta_1 \quad \dots \quad \theta_n)^T$  and let the  $i$ -th data vector be:  $x^{(i)} = (x_0 \quad x_1 \quad \dots \quad x_n)^T$  where  $x_0 = 1$ .

### 1.1 Write the hypothesis function $h_\theta(x)$ as a vectorial expression.

The summation notation for the hypothesis function is:

$$h_\theta(x^{(i)}) = \theta_0 x_0^{(i)} + \theta_1 x_1^{(i)} + \theta_2 x_2^{(i)} + \dots + \theta_n x_n^{(i)}$$

This is the same as the result of the following matrix multiplication:

$$h_\theta(x^{(i)}) = \theta^T x^{(i)}, \tag{1}$$

which is a vectorial expression.

### 1.2 What is the vectorized expression for the cost function: $J(\theta)$ ?

The cost function in the notation used up to now:

$$J(\theta) = \frac{1}{2m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})^2.$$

If we simply insert the vectorial notation of the hypothesis function from last question (1) we get:

$$J(\theta) = \frac{1}{2m} \sum_{i=1}^m (\theta^T x^{(i)} - y^{(i)})^2.$$

### 1.3 What is the vectorized expression for the gradient of the cost function?

*i.e. what is:*

$$\frac{\delta J(\theta)}{\delta \theta} = \begin{pmatrix} \frac{\delta J(\theta)}{\delta \theta_0} \\ \vdots \\ \frac{\delta J(\theta)}{\delta \theta_n} \end{pmatrix}$$

Again the explicit summation over the data vectors from the learning set is allowed here.

The notation we used up until now is:

$$\frac{\delta J(\theta)}{\delta \theta_j} = \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_j^{(i)}.$$

If we integrate the vectorized notation of the hypothesis (1) again we get:

$$\frac{\delta J(\theta)}{\delta \theta_j} = \frac{1}{m} \sum_{i=1}^m (\theta^T x^{(i)} - y^{(i)}) x_j^{(i)}.$$

Substituting this summation for the  $\delta$  notation will give us the following vector:

$$\frac{\delta J(\theta)}{\delta \theta} = \frac{1}{m} \begin{pmatrix} \sum_{i=1}^m (\theta^T x^{(i)} - y^{(i)}) \\ \sum_{i=1}^m (\theta^T x^{(i)} - y^{(i)}) x_1^{(i)} \\ \vdots \\ \sum_{i=1}^m (\theta^T x^{(i)} - y^{(i)}) x_n^{(i)} \end{pmatrix} \quad (2)$$

### 1.4 What is the vectorized expression for the $\theta$ update rule in the gradient descent procedure?

The original notation for the update rule for one theta was:

$$\theta_j = \theta_j - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_i^{(i)}$$

By writing it in vector notation we update the entire theta instead of each element separately. Using the formulas from (1) and (2) we get:

$$\theta = \theta - \alpha \frac{1}{m} \begin{pmatrix} \sum_{i=1}^m (\theta^T x^{(i)} - y^{(i)}) \\ \sum_{i=1}^m (\theta^T x^{(i)} - y^{(i)}) x_1^{(i)} \\ \vdots \\ \sum_{i=1}^m (\theta^T x^{(i)} - y^{(i)}) x_n^{(i)} \end{pmatrix}$$

### 1.5 (BONUS) Remove the explicit summation by using a matrix vector multiplication

We start by defining the data matrix  $X$ ; every row of  $X$  is a training example, the first column containing  $x_0^{(1)}, x_0^{(2)}, \dots, x_0^{(n)}$ .

$$X = \begin{pmatrix} x_0^{(1)} & x_1^{(1)} & \dots & x_n^{(1)} \\ x_0^{(2)} & x_1^{(2)} & \dots & x_n^{(2)} \\ \vdots & \vdots & & \vdots \\ x_0^{(m)} & x_1^{(m)} & \dots & x_n^{(m)} \end{pmatrix}$$

The hypothesis function will then be  $h_\theta(X) = X\theta$ , which results in the matrix:

$$h_\theta(X) = X\theta = \begin{pmatrix} x_0^{(1)} & x_1^{(1)} & \dots & x_n^{(1)} \\ x_0^{(2)} & x_1^{(2)} & \dots & x_n^{(2)} \\ \vdots & \vdots & & \vdots \\ x_0^{(m)} & x_1^{(m)} & \dots & x_n^{(m)} \end{pmatrix} \begin{pmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_n \end{pmatrix} = \begin{pmatrix} \theta_0 x_0^{(1)} + \theta_1 x_1^{(1)} + \dots + \theta_n x_n^{(1)} \\ \theta_0 x_0^{(2)} + \theta_1 x_1^{(2)} + \dots + \theta_n x_n^{(2)} \\ \vdots \\ \theta_0 x_0^{(m)} + \theta_1 x_1^{(m)} + \dots + \theta_n x_n^{(m)} \end{pmatrix}$$

Let  $Y$  be the matrix  $Y = (y^{(1)} \ y^{(2)} \ \dots \ y^{(m)})^T$ . Then we can write the derivative of  $J(\theta)$  as follows:

$$\begin{aligned} \frac{\delta J(\theta)}{\delta \theta} &= \frac{1}{m} X(X\theta - Y) \\ &= \frac{1}{m} \begin{pmatrix} x_0^{(1)} & x_1^{(1)} & \dots & x_n^{(1)} \\ x_0^{(2)} & x_1^{(2)} & \dots & x_n^{(2)} \\ \vdots & \vdots & & \vdots \\ x_0^{(m)} & x_1^{(m)} & \dots & x_n^{(m)} \end{pmatrix} \left( \begin{pmatrix} \theta_0 x_0^{(1)} + \theta_1 x_1^{(1)} + \dots + \theta_n x_n^{(1)} \\ \theta_0 x_0^{(2)} + \theta_1 x_1^{(2)} + \dots + \theta_n x_n^{(2)} \\ \vdots \\ \theta_0 x_0^{(m)} + \theta_1 x_1^{(m)} + \dots + \theta_n x_n^{(m)} \end{pmatrix} - \begin{pmatrix} y^{(1)} \\ y^{(2)} \\ \vdots \\ y^{(m)} \end{pmatrix} \right) \end{aligned}$$

Thus, the final update rule would be:

$$\theta = \theta - \alpha \frac{1}{m} X(X\theta - Y)$$