

1 Summary

Let's start from linear regression with L2 regularization [1].

$$\mathbf{y} = \mathbf{w}^T \phi(\mathbf{x})$$

Note that $\phi(\mathbf{x}) = X$ where X is a $m \times D$ matrix where each training sample has m features.

To choose the optimal \mathbf{w} , we minimize the sum of squares:

$$\mathbf{w}^* = \arg \min_{\mathbf{w}} \|\mathbf{y} - \mathbf{w}^T \phi(\mathbf{x})\|^2$$

To avoid overfitting, we add the regularization parameter:

$$\mathbf{w}^* = \arg \min_{\mathbf{w}} \|\mathbf{y} - \mathbf{w}^T \phi(\mathbf{x})\|^2 + \lambda \|\mathbf{w}\|^2$$

Lets say that $L(\mathbf{w}) = \|\mathbf{y} - \mathbf{w}^T \phi(\mathbf{x})\|^2 + \lambda \|\mathbf{w}\|^2$

What is this derivative? We will compute the following gradient:

$$\nabla_{\mathbf{w}} L(\mathbf{w}) = \frac{\partial L(\mathbf{w})}{\partial \mathbf{w}} = \left(\frac{\partial L(\mathbf{w})}{\partial w_1}, \dots, \frac{\partial L(\mathbf{w})}{\partial w_i} \right)$$

Let's go look into one element of this gradient $\nabla_{\mathbf{w}}$:

$$\frac{\partial L(\mathbf{w})}{\partial w_i} = \frac{\partial (\|\mathbf{y} - \mathbf{w}^T \phi(\mathbf{x})\|^2 + \lambda \|\mathbf{w}\|^2)}{\partial w_i}$$

Keep in mind that

$$\|\mathbf{y} - X\mathbf{w}\|^2 = \mathbf{y}^T \mathbf{y} - \mathbf{y}^T X\mathbf{w} - (X\mathbf{w})^T \mathbf{y} + (X\mathbf{w})^T (X\mathbf{w}) \quad (1)$$

$$= \mathbf{y}^T \mathbf{y} - 2\mathbf{y}^T X\mathbf{w} + (X\mathbf{w})^T (X\mathbf{w}) \quad (2)$$

$$= \sum_i y_i^2 - 2\mathbf{y}^T X\mathbf{w} + \sum_j (\mathbf{x}_j \mathbf{w})^2 \quad (3)$$

$$= \sum_i y_i^2 - 2 \sum_i a_i w_i + \sum_j \left(\sum_i x_{ji} w_i \right)^2 \quad (4)$$

$$= \sum_i y_i^2 - 2 \sum_i a_i w_i + \sum_j (x_{j1} w_1 + \dots + x_{ji} w_i)^2 \quad (5)$$

where $\mathbf{a} = \mathbf{y}^T X$, $a_i = \mathbf{y}^T X_i$

Since $\forall k$ s.t. $k \neq i$, $\frac{\partial w_k}{\partial w_i} = 0$,

$$\frac{\partial L(\mathbf{w})}{\partial w_i} = -2a_i - \frac{\partial \sum_j ((x_{ji}w_i)(x_{j1}w_1 + \dots + x_{ji-1}w_{i-1}) + (x_{j1}w_1 + \dots + x_{ji-1}w_{i-1})(x_{ji}w_i) + (x_{ji}w_i)^2)}{\partial w_i} \quad (6)$$

$$= -2a_i - \sum_j ((x_{ji})(x_{j1}w_1 + \dots + x_{ji-1}w_{i-1}) + (x_{j1}w_1 + \dots + x_{ji-1}w_{i-1})(x_{ji}) + (2x_{ji}^2w_i)) \quad (7)$$

$$= -2a_i - \sum_j 2((x_{ji})(x_{j1}w_1 + \dots + x_{ji-1}w_{i-1}) + (2x_{ji}^2w_i)) \quad (8)$$

$$= -2(a_i - \sum_j ((x_{ji})(x_{j1}w_1 + \dots + x_{ji-1}w_{i-1}) + (x_{ji}^2w_i))) \quad (9)$$

$$= -2(a_i - \sum_j ((x_{ji})(x_{j1}w_1 + \dots + x_{ji-1}w_{i-1} + x_{ji}w_i))) \quad (10)$$

$$= -2(a_i - \sum_j ((x_{ji})(x_{j1}, \dots, x_{ji-1}, x_{ji})\mathbf{w})) \quad (11)$$

$$= -2(\mathbf{y}^T X_i - X_j^T X_i \mathbf{w}) \quad (12)$$

In the matrix representation:

$$\nabla_{\mathbf{w}} L(\mathbf{w}) = -2(\mathbf{y}^T X - X^T X \mathbf{w})$$

Also in general:

$$\nabla_{\mathbf{w}} \mathbf{w}^T X^T X \mathbf{w} = 2X^T X \mathbf{w}$$

Note that $X^T X$ is symmetric and it is part of the common matrix derivative pattern¹

References

- [1] Lecture 11: Regularization. <http://grandmaster.colorado.edu/~cketelsen/files/csci5622/videos/lesson11/lesson11.pdf>. [Online; accessed 19-Nov-2016].

¹https://en.wikipedia.org/wiki/Matrix_calculus#Scalar-by-vector_identities