

Ridge regression

Nathanaël Carraz Rakotonirina

Mathématiques Informatique et Statistique Appliquées (MISA)
Université d'Antananarivo

To avoid overfitting, a common strategy is to constrain the weights to be smaller. Weights that become too large in magnitude are penalized. This is called **weight decay** or L^2 **regularization**.

The loss function is updated as follows

$$J(\mathbf{w}) = \text{RSS}(\mathbf{w}) + \lambda \|\mathbf{w}\|_2^2$$

where λ controls the strength of the regularizer

Probabilistic interpretation

Instead of MLE, we use MAP estimation with a zero-mean Gaussian prior on the weights $p(\mathbf{w}) = \prod_{i=1}^D \mathcal{N}(w_i; 0, \tau^2)$. The MAP estimate corresponds to minimizing

$$J(\mathbf{w}) = \frac{1}{2\sigma^2} \|\mathbf{X}\mathbf{w} - \mathbf{y}\|_2^2 + \frac{1}{2\tau^2} \|\mathbf{w}\|_2^2$$

where $\lambda = \frac{\sigma^2}{\tau^2}$

$$J(\mathbf{w}) = \|\mathbf{X}\mathbf{w} - \mathbf{y}\|_2^2 + \lambda \|\mathbf{w}\|_2^2 = \text{RSS}(\mathbf{w}) + \lambda \|\mathbf{w}\|_2^2$$

Solving the MAP estimate

The gradient is given by

$$\nabla_{\mathbf{w}} J(\mathbf{w}) = 2(\mathbf{X}^\top \mathbf{X} \mathbf{w} - \mathbf{X}^\top \mathbf{y} + \lambda \mathbf{w})$$

we set the gradient to zero to obtain

$$\hat{\mathbf{w}} = (\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I}_D)^{-1} \mathbf{X}^\top \mathbf{y}$$

As in standard linear regression, computing directly the inverse is slow and could be numerically unstable. There are other ways.

- ▶ Convert it to standard least squares then use QR decomposition
- ▶ SVD