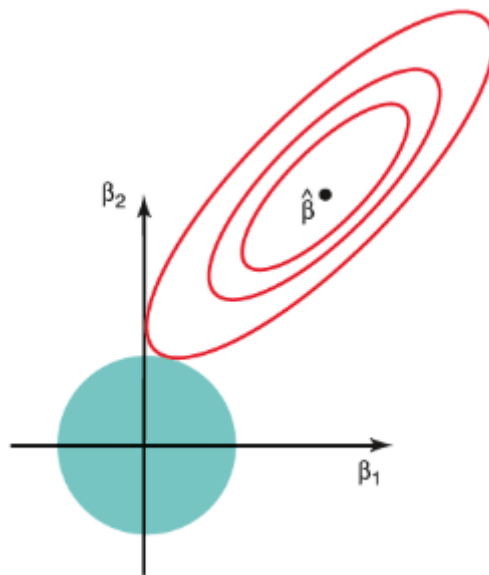# RAHUL JAIN

Follow        14 Followers        About

# L2 Regularisation: Maths

RAHUL JAIN  Feb 8, 2020 · 3 min read

L2 regularisation, one of the most popular techniques in Machine Learning, is a method to reduce the variance of the model and increase the bias to make the model more generalisable.



Geometrical Representation of L2 Regularisation

L2 regularisation is a MAP estimation with Gaussian prior probabilities. Due to the imposed prior, the model can generalise the data well by scaling the weights as per their significance.

The L2 regularised objective function and its gradient is given by,

$$\nabla_w \tilde{J}(w; X, y) = \alpha w + \nabla_w J(w; X, y).$$

From Taylor series expansion, we will make a quadratic approximation of the unregularized objective function around the point $w^*$. $w^*$ is the minimiser of the unregularized objective function.

$$\hat{J}(\theta) = J(w^*) + \frac{1}{2}(w - w^*)^\top H(w - w^*),$$

$$\text{Since} \quad w^* = \arg\min_w J(w)$$

$$\nabla_w J(w; X, y)\Big|_{w^*} = 0$$

Thus, the regularised objective function and its gradient becomes,

$$\tilde{J}(w) = \frac{\alpha}{2} w^\top w + J(w^*) + \frac{1}{2}(w - w^*)^\top H(w - w^*),$$

$$\nabla_w \tilde{J}(w) = \alpha w + H(w - w^*)$$

Let $\tilde{w}$ represents the minimiser of the regularised objective function. The slope at this point is equal to zero.

$$\nabla_w \tilde{J}(w) = 0$$

$$\alpha \tilde{w} + H(\tilde{w} - w^*) = 0$$

$$(H + \alpha I)\tilde{w} = H w^*$$

$$\tilde{w} = (H + \alpha I)^{-1} H w^*.$$

Since H is positive semi-definite, real and symmetric matrix, it can be decomposed into a diagonal matrix $\Lambda$ and an orthonormal basis of eigenvectors, Q, such that,

$$\tilde{w} = (Q\Lambda Q^\top + \alpha I)^{-1} Q\Lambda Q^\top w^*$$
$$= \left[ Q(\Lambda + \alpha I)Q^\top \right]^{-1} Q\Lambda Q^\top w^*$$
$$= Q(\Lambda + \alpha I)^{-1}\Lambda Q^\top w^*.$$

The effect of weight decay is to rescale $w*$ along the axes defined by the eigenvectors of H. The component of $w\sim$ along the ith eigenvector of H is rescaled.

$$\tilde{w}_i = \frac{\lambda_i}{\lambda_i + \alpha} w^*_i$$

If $\lambda i >> \alpha$, the effect of regularisation is relatively small. However, components with $\lambda i << \alpha$ will be shrunk to have nearly zero magnitudes. Only directions along which the parameters contribute significantly to reducing the objective function are preserved relatively intact. In other unimportant instructions, indicated by a small eigenvalue of the Hessian, weight vectors are decayed away through the use of the regularisation throughout the training.

The regularisation constant, $\alpha$ is a hyperparameter and is tuned to get the best results. As the value of $\alpha$ increases, weights are decayed more.

## References:

Deep Learning (Adaptive Computation and Machine Learning series) — By Ian Goodfellow, Yoshua Bengio, Aaron Courville

L2 Regularization      Machine Learning      Maths      Gradient Descent      Proof