

3) Elastic Net

is a combination of L_2 and L_1 .

The implementation has both λ_1 and λ_2 for L_1 and L_2 respectively some $\lambda_1 = \lambda_2$

$$L(\lambda_1, \lambda_2, \beta) = |y - X\beta|^2 + \lambda_2 |\beta|^2 + \lambda_1 |\beta|_1$$

In logistic regression the $|y - X\beta|^2$ is $J(\beta)$

$$J(\beta) = -y_i \log h(\beta) - (1 - y_i) \log (1 - h(\beta))$$

β is the weight for features we need to learn.

So

$$L(\lambda_1, \lambda_2, \beta) = \sum_{i=1}^n -y_i \log h(\beta) - (1 - y_i) \log (1 - h(\beta)) + \lambda_2 |\beta|^2 + \lambda_1 |\beta|_1$$

Where $\frac{L(\lambda_1, \lambda_2, \beta)}{d\beta} \Rightarrow$ to find the maximum solution for gradient descent iterative approach for β

$$|\beta|^2 = \sum_{k=1}^d \beta_k^2$$

$$|\beta|_1 = \sum_{k=1}^d |\beta_k|$$

d is number of features.

(1)

Gradient descent step:

$$\frac{\partial L}{\partial w_{ij}} = \sum_{i=1}^n x_{ij} (y_i - \hat{f}(\beta)) + \lambda_2 |\beta| + \lambda_1 \frac{|\beta|}{\beta}$$

The gradient descent step is divided by samples to average out.

The L_2 is applied with the data and after that L_1 is applied on the previous weights without updating ~~the~~ weights learned from L_2 . The elastic net estimator is a 2-stage procedure: for fixed λ we first find ridge regression coefficients (β_L) and then we do the lasso shrinkage along the lasso coefficient solution path.