

Working on Real Big Data

Regularization

Since the number of features is far larger than the number of observation, to prevent overfitting, regularization restricts the freedom of a model's parameters, penalizing their distance from some prior belief. Widely used regularizers penalize large weights with an objective function of the form

$$F(\beta) = L(\beta) + R(\beta) \quad (1)$$

where β is the vector of coefficients of interested, and $L(\beta)$ is the loss function with respect to β . Many commonly used regularizers $R(\beta)$ are of the form $\lambda \|\beta\|$ where λ determines the strength of regularization and the ℓ_0 , ℓ_1 , ℓ_2^2 , or ℓ_∞ are common choices for the penalty function.

Lazy Update

For a particular observation, the majority of feature value is zero. To address sparsity, lazy update is applied. The idea is to update a coefficient only when the feature is nonzero. The essence of the approach is given as follows.

Algorithm 1 Lazy Updates

Require: $\psi \in \mathbb{R}^d$
for $t \in 1, \dots, T$ **do**
 Sample x_i randomly from the training set
 for j s.t. $x_{ij} \neq 0$ **do**
 $w_j \leftarrow \text{Lazy}(w_j, t, \psi_j)$
 $\psi_j \leftarrow t$
 end for
 $w \leftarrow w - \nabla F_i(w)$
end for

We keep an array $\psi \in \mathbb{R}^d$ in which each ψ_j stores the index of the last iteration at which the value of weight j was updated. When processing observation x_i at iteration k , we iterate through its nonzero features x_{ij} . For each such nonzero feature, we lazily apply the $k - \psi_j$ delayed updates collectively. Using the updated weights, we compute the prediction $y^{(k)}$ with the fully updated relevant parameters from $\beta^{(k)}$. We then compute the gradient and update these parameters. For **Ridge SGD Sparse Weight Update**, we have,

$$\beta_t^i = \beta_{t-n}^i (1 - \alpha\lambda)^n \quad (2)$$

where t is the current iteration, n is the collective steps we need to update, which is the difference between t and the last time β^i was updated, α is learning rate, and λ is the strength of regularization. For **Lasso SGD Sparse Weight Update**, we have

$$\beta_t^i = \beta_{t-n}^i \max(0, 1 - \frac{n\alpha\lambda}{|\beta_{t-n}^i|}) \quad (3)$$