

Cross-entropy loss with regularization and its gradients

Notation:

n - number of train examples (batch_size);

d - number of dimensions of an example;

c - number of classes (10 classes for CIFAR-10);

y - labels of correct classes (y_i - the correct label for the i th example);

$$L = \frac{1}{n} \sum_i^n L_i + \lambda * \sum_j^d \sum_k^c W_{j,k}^2$$

$$L_i = -\ln p_{y_i} [\text{cross-entropy}]$$

$$p_k = \frac{e^{z_k}}{\sum_j^c e^{z_j}} [\text{softmax}]$$

$$z_j = X * W_j [\text{class-score} - z_j - \text{for-example} - X]$$

$$\frac{\partial L_i}{\partial p_{y_i}} = -\frac{1}{p_{y_i}}$$

$$\frac{\partial p_k}{\partial z_k} = \frac{e^{z_k} \sum_j^c e^{z_j} - e^{z_k} e^{z_k}}{(\sum_j^c e^{z_j})^2} = \frac{e^{z_{y_i}}}{\sum_j^c e^{z_j}} - \frac{e^{z_k}}{\sum_j^c e^{z_j}} \frac{e^{z_k}}{\sum_j^c e^{z_j}} = p_k - p_k^2 = p_k(1 - p_k)$$

$$\frac{\partial p_k}{\partial z_l} = -\frac{e^{z_k} e^{z_l}}{(\sum_j^c e^{z_j})^2} = -\frac{e^{z_k}}{\sum_j^c e^{z_j}} \frac{e^{z_l}}{\sum_j^c e^{z_j}} = -p_k p_l [l \neq k]$$

$$\frac{\partial z_j}{\partial w_j} = X$$

$$\frac{\partial L_i}{\partial w_{y_i}} = \frac{\partial L_i}{\partial p_{y_i}} \frac{\partial p_{y_i}}{\partial z_{y_i}} \frac{\partial z_{y_i}}{\partial w_{y_i}} = -\frac{1}{p_{y_i}} p_{y_i} (1 - p_{y_i}) X = (p_{y_i} - 1) X [y_i, \text{chain-rule}]$$

$$\frac{\partial L_i}{\partial w_j} = \frac{\partial L_i}{\partial p_{y_i}} \frac{\partial p_{y_i}}{\partial z_j} \frac{\partial z_j}{\partial w_j} = -\frac{1}{p_{y_i}} (-p_{y_i} p_j) X = p_j X [j \neq y_i, \text{chain-rule}]$$