



## Module 2

### Adagrad

Adagrad <sup>[9]</sup> is an algorithm for gradient-based optimization that does just this: It adapts the learning rate to the parameters, performing smaller updates (i.e. low learning rates) for parameters associated with frequently occurring features, and larger updates (i.e. high learning rates) for parameters associated with infrequent features. For this reason, it is well-suited for dealing with sparse data. Dean et al. <sup>[10]</sup> have found that Adagrad greatly improved the robustness of SGD and used it for training large-scale neural nets at Google, which -- among other things -- learned to recognize cats in Youtube videos. Moreover, Pennington et al. <sup>[11]</sup> used Adagrad to train GloVe word embeddings, as infrequent words require much larger updates than frequent ones.

Previously, we performed an update for all parameters  $\theta$  at once as every parameter  $\theta_i$  used the same learning rate  $\eta$ . As Adagrad uses a different learning rate for every parameter  $\theta_i$  at every time step  $t$ , we first show Adagrad's per-parameter update, which we then vectorize. For brevity, we use  $g_t$  to denote the gradient at time step  $t$ .  $g_{t,i}$  is then the partial derivative of the objective function w.r.t. to the parameter  $\theta_i$  at time step  $t$ :

$$g_{t,i} = \nabla_{\theta} J(\theta_{t,i}).$$

The SGD update for every parameter  $\theta_i$  at each time step  $t$  then becomes:

$$\theta_{t+1,i} = \theta_{t,i} - \eta \cdot g_{t,i}.$$

In its update rule, Adagrad modifies the general learning rate  $\eta$  at each time step  $t$  for every parameter  $\theta_i$  based on the past gradients that have been computed for  $\theta_i$ :



## Module 2

$$\theta_{t+1,i} = \theta_{t,i} - \frac{\eta}{\sqrt{G_{t,ii}} + \epsilon} \cdot g_{t,i}.$$

$G_t \in \mathbb{R}^{d \times d}$  here is a diagonal matrix where each diagonal element  $i$ ,  $i$  is the sum of the squares of the gradients w.r.t.  $\theta_i$  up to time step  $t$  <sup>[12]</sup>, while  $\epsilon$  is a smoothing term that avoids division by zero (usually on the order of  $1e - 8$ ). Interestingly, without the square root operation, the algorithm performs much worse.

As  $G_t$  contains the sum of the squares of the past gradients w.r.t. to all parameters  $\theta$  along its diagonal, we can now vectorize our implementation by performing a matrix-vector product  $\odot$  between  $G_t$  and  $g_t$ :

$$\theta_{t+1} = \theta_t - \frac{\eta}{\sqrt{G_t} + \epsilon} \odot g_t.$$

One of Adagrad's main benefits is that it eliminates the need to manually tune the learning rate. Most implementations use a default value of 0.01 and leave it at that.

Adagrad's main weakness is its accumulation of the squared gradients in the denominator: Since every added term is positive, the accumulated sum keeps growing during training. This in turn causes the learning rate to shrink and eventually become infinitesimally small, at which point the algorithm is no longer able to acquire additional knowledge. The following algorithms aim to resolve this flaw.