CS 6.5: Theories of Deep Learning Problem Sheet 3

Prof. Jared Tanner November 19, 2019

Gradient-Descent based optimization for neural networks

Gradient descent is one of the most popular algorithms to perform optimization and by far the most common way to optimize neural networks. Gradient descent is a way to minimize an objective function $J(\theta)$ parametrized by a model's parameters θ by updating the parameters in the opposite direction of the gradient of the objective function $\nabla J(\theta)$ w.r.t. to the parameters. The learning rate η determines the size of the steps we take to reach a (local) minimum. In other words, we follow the direction of the slope of the surface created by the objective function downhill until we reach a valley.

The aims of this exercise is to provide you with intuitions towards the behaviour of such optimizing algorithms and associated challenges during actual training in practice. While, in class we will discuss a few of such algorithms, this exercise focuses on 'ADADELTA'. ADADELTA, try to adjust the learning rate during training by annealing, i.e. reducing the learning rate according to a pre-defined criteria. However, the same learning rate doesn't applies to all the dimensions of θ , and depends on the frequency of their use in overall learning process. Additionally, with ADADELTA, we do not even need to set a default learning rate η , as it has been eliminated from the update rule.

1. Task1: Weight Decay: One of the ways to update the parameters θ is to consider the updates of form:

$$\theta_{t+1} = (1 - \lambda)\theta_t - \eta \nabla J_t(\theta_t), \tag{1}$$

where λ defines the rate of the weight decay per step and $\nabla J_t(\theta_t)$ is the t-th batch gradient to be multiplied by a learning rate η . Prove that the standard SGD with base learning rate η executes the same steps on $J_t(\theta)$ with weight decay λ (defined in Equation 1) as it executes without weight decay on $J_t'(\theta) = J_t(\theta) + \frac{\lambda'}{2} \|\theta\|_2^2$, with $\lambda' = \frac{\lambda}{\eta}$.

- 2. Task2: Weight Decay for adaptive gradients based approaches² (such as ADADELTA, ADAGRAD, ADAM): Let O be an optimizer that iterates as $\theta_{t+1} = \theta_t - \eta \mathbf{M}_t J_t(\theta_t)$ without weight decay, and $\theta_{t+1} = (1 - \lambda)\theta_t - \eta \mathbf{M}_t J_t(\theta_t)$ with weight decay, with $\mathbf{M}_t \neq k\mathbf{I}, k \in \mathbb{R}$. Prove that for O there exists no coefficient λ such that running O on $J_t'(\theta) = J_t(\theta) + \frac{\lambda'}{2} ||\theta||_2^2$ without weight decay is equivalent to running O on $J_t(\theta)$ with decay $\lambda \in \mathbb{R}^+$
- 3. Task3: Weight Decay and Batch-Normalization (BN).

The motivation for Weight Decay is to reduce overfitting and for BN is to stabilize the distribution of layer activations over the course of training. The key property of BN is that it make the neural net output approximately invariant to the scale of previous layer output. This brings us to the key question: What Happens When Both Are Used Together?

Show that BN cancels the effect of Weight Decay as it no longer can directly alter any dimension of weight vector in future gradient descent steps. Also discuss the impact on effective learning rate in this scenario. (Note: non-rigorous derivations are fine for this task)

¹ADADELTA:https://arxiv.org/pdf/1212.5701.pdf

²Survey on GD:https://arxiv.org/pdf/1609.04747.pdf