

Series Monday, Nov 5, 2018 (Deep Learning, Exercise series 6)

Problem 1 (Fundamentals of Unconstrained Optimization):

1. Compute the gradient $\nabla f(w)$ and Hessian $\nabla^2 f(w)$ of the Rosenbrock function

$$f(w) = 100(w_2 - w_1^2)^2 + (1 - w_1)^2.$$

Show that $w^* = (1, 1)^\top$ is the only local minimizer of this function, and that the Hessian matrix at that point is positive definite.

2. Show that the function $f(w) = 8w_1 + 12w_2 + w_1^2 - 2w_2^2$ has only one stationary point, and that it is neither a maximum or minimum, but a saddle point. Sketch the contour lines of f .

Problem 2 (Gradient Descent step derivation):

Show that the Gradient Descent step in iteration t , namely

$$s_t := w_{t+1} - w_t = -\frac{1}{L} \nabla f(w_t), \quad L > 0 \quad (1)$$

minimizes a quadratically regularized first order Taylor expansion of the objective f around the current iterate w_t . That is, show that the step s_t^* defined as

$$s_t^* := \arg \min_{s \in \mathbb{R}^d} \left[m_t(s) := f(w_t) + s^\top \nabla f(w_t) + \frac{L}{2} \|s\|_2^2 \right] \quad (2)$$

equals s_t as given in Eq. (1). (As we will see next week, if we pick L greater or equal to the Lipschitz constant of the gradients, $m_t(s)$ is a global overestimator of f and hence s_t is guaranteed to be a descent step.)

Problem 3 (Programming exercise: Optimization methods in tensorflow):

In this exercise, you will implement three of the most common optimization methods in deep learning—SGD, Adagrad, and Adam—from scratch. In the supplementary material of this exercise, we provide a jupyter notebook skeleton for a simple MNIST classification model. The training and testing procedure is already implemented. The only thing left to do is the implementation of an optimization method. But instead of using the tensorflow API (e.g. `tf.train.AdamOptimizer()`) you are asked to implement the optimizers yourself. A skeleton of the necessary classes and methods is provided.

1. Implement Stochastic Gradient Descent (**SGD**) in the corresponding cell. The SGD update step is given by

$$w_{t+1} = w_t - \eta \nabla_{w_t} l(x_t; w_t) \quad (3)$$

where $w \in \mathbb{R}^n$ are the model parameters, x_t is a randomly sampled *mini-batch* of the data and $l()$ your loss function. The provided signature of the method `apply_gradients` follows the naming convention for tensorflow optimizers. Try to first get familiar with how the gradients and variables are provided and how to assign new values to tensorflow variables.

2. Implement **Adagrad** [1] in the corresponding cell. Its update step is computed by

$$g_t = \nabla_w l(x_t; w_t) \quad (4)$$

$$G_t = G_{t-1} + g_t g_t^T \in \mathbb{R}^{n \times n} \quad (5)$$

$$w_{t+1} = w_t - \eta \text{diag}(G_t + \epsilon I)^{-0.5} g_t \quad (6)$$

where $\mathbf{w} \in \mathbb{R}^n$ are the model parameters, \mathbf{x}_t is a randomly sampled *mini-batch* of the data and $l()$ your loss function. As shown by the formulas, Adagrad uses individual step sizes for every parameter, which depend on the history of its gradients. The intuition here is that the learning rate decays faster for dimensions that have already seen significant updates.

Hint: For the implementation it is not necessary to compute the full matrix G_t but just its diagonal elements, which reduces to the sum of squared gradients in every dimension over time.

3. Implement **Adam** [2] in the corresponding cell. Its update step¹ is computed by

$$\mathbf{m}_t = \beta_1 \mathbf{m}_{t-1} + (1 - \beta_1) \nabla_{\mathbf{w}} l(\mathbf{x}_t; \mathbf{w}_t) \quad (7)$$

$$\mathbf{v}_t = \beta_2 \mathbf{v}_{t-1} + (1 - \beta_2) (\nabla_{\mathbf{w}} l(\mathbf{x}_t; \mathbf{w}_t))^2 \quad (8)$$

$$\eta_t = \eta * \frac{\sqrt{1 - \beta_2^t}}{1 - \beta_1^t} \quad (9)$$

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \frac{\eta_t}{\sqrt{\mathbf{v}_t} + \epsilon} \mathbf{m}_t \quad (10)$$

where $\mathbf{w} \in \mathbb{R}^n$ are the model parameters, \mathbf{x}_t is a randomly sampled *mini-batch* of the data and $l()$ your loss function.

References

- [1] John Duchi, Elad Hazan, and Yoram Singer. Adaptive subgradient methods for online learning and stochastic optimization. Technical Report UCB/EECS-2010-24, EECS Department, University of California, Berkeley, Mar 2010.
- [2] Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *CoRR*, abs/1412.6980, 2014.

¹Note that the following formula corresponds to Algorithm 1 in [2] with the proposed modification of section 2. This is the same algorithm that is used in the tensorflow implementation.