

## OPTIMIZATION. HOMEWORK 12

OSCAR DALMAU

(1) Consider the optimization problem

$$F(\theta) = \frac{1}{N} \sum_{i=1}^N (h_{\theta}(\mathbf{x}_i) - y_i)^2 \quad (1)$$

where  $(\mathbf{x}_i, y_i)$ ,  $\mathbf{x}_i \in \mathbb{R}^n$ ,  $y_i \in \{0, 1\}$ ,  $i = 1, 2, \dots, N$  are given and

$$h_{\theta}(\mathbf{x}) = f_{\mathbf{a},b}(g_{\mathbf{c},\mathbf{d}}(\mathbf{x}))$$

$$g_{\mathbf{c},\mathbf{d}} : \mathbb{R}^n \rightarrow \mathbb{R}^m$$

$$f_{\mathbf{a},b} : \mathbb{R}^m \rightarrow \mathbb{R}$$

$m \in \mathbb{N}$  is known,

$$g_{\mathbf{c},\mathbf{d}}(\mathbf{x}) = [\sigma(\mathbf{c}_j^T \mathbf{x} + d_j)]_{j=1}^m$$

$$f_{\mathbf{a},b}(\mathbf{z}) = \sigma(\mathbf{a}^T \mathbf{z} + b)$$

$$\sigma(t) = \frac{1}{1 + e^{-t}}, \quad t \in \mathbb{R}$$

and  $\theta$  corresponds to the set of parameters  $\mathbf{a}, b, \mathbf{c}, \mathbf{d}$ .

Implement the following algorithms:

- Gradient descent with fixed step size and gradient approximation.
- BFGS with gradient approximation.

Compare the previous algorithms with respect to number of iteration and computational time using the mnist dataset and selecting only two digits.

**Note:** You can use any library to load mnist dataset. However, you can use dataset mnist.pkl.gz

- `train_set[0]` is a matrix of size  $(50000, 784)$  where  $n = 50000$  is the number of observations and each row represents an observation  $\mathbf{x}_i \in \mathbb{R}^{784}$

- `train_set[1]` is a vector of size (50000) where each entry  $y_i \in \{0, 1, \dots, 9\}$

Select from `train_set[0]` and `train_set[1]` the set of observations  $\mathcal{S} = \{(\mathbf{x}_i, y_i)\}$  with  $\mathbf{x}_i \in \mathbb{R}^{784}$  and  $y_i \in \{0, 1\}$  and estimate the parameters  $\hat{\theta}$  that minimizes the function  $F(\theta)$  in equation (1).

Select from `test_set[0]` and `test_set[1]` the set  $\mathcal{T} = \{(\mathbf{x}_i, y_i)\}$  such that  $\mathbf{x}_i \in \mathbb{R}^{784}$  and  $y_i \in \{0, 1\}$  and compute the error

$$error = \frac{1}{|\mathcal{T}|} \sum_{(\mathbf{x}_i, y_i) \in \mathcal{T}} |\mathbf{1}_{h_{\hat{\theta}}(\mathbf{x}_i) > 0.5}(\mathbf{x}_i) - y_i|$$

where  $|\mathcal{T}|$  represents the number of elements of the set  $\mathcal{T}$ .