

CS 541 Artificial Intelligence: Homework 3

Instructor: Jie Shen

Due: 11/08/2020, 20:00 pm EST

Gradient Calculation

Suppose $\mathbf{x} \in \mathbb{R}^d$ and $y \in \mathbb{R}$ are known. Calculate the gradient of the following functions:

- Sigmoid function $F(\mathbf{w}) = 1/(1 + e^{-\mathbf{x} \cdot \mathbf{w}})$;
- Logistic loss $F(\mathbf{w}) = \log(1 + e^{-y\mathbf{x} \cdot \mathbf{w}})$.

Note: we take both \mathbf{x} and \mathbf{w} as column vectors.

Linear Regression

Suppose we are given a data set $\{\mathbf{x}_i, y_i\}_{i=1}^n$, where each $\mathbf{x}_i \in \mathbb{R}^d \times \mathbb{R}$ is a row vector. We hope to learn a mapping f such that each y_i is approximated by $f(\mathbf{x}_i)$. Then a popular approach is to fit the data with *linear regression* – it assumes there exists $\mathbf{w} \in \mathbb{R}^d$ such that $y_i \approx \mathbf{w} \cdot \mathbf{x}_i$. In order to learn \mathbf{w} from the data, it typically boils down to solving the following *least-squares* program:

$$\min_{\mathbf{w} \in \mathbb{R}^d} F(\mathbf{w}) := \frac{1}{2} \|\mathbf{y} - \mathbf{X}\mathbf{w}\|_2^2, \quad (1)$$

where \mathbf{X} is the data matrix with the i th row being \mathbf{x}_i , and $\mathbf{y} = (y_1, y_2, \dots, y_n)^\top$.

1. Compute the gradient and the Hessian matrix of $F(\mathbf{w})$, and show that (1) is a convex program.
2. Note that (1) is equivalent to the following:

$$\min_{\mathbf{w} \in \mathbb{R}^d} \|\mathbf{y} - \mathbf{X}\mathbf{w}\|_2^{100},$$

in the sense that any minimizer of (1) is also an optimum of the above, and vice versa. State why we stick with the least-squares formulation.

3. State a possible condition on \mathbf{X} such that $F(\mathbf{w})$ is strongly-convex. Under which condition on \mathbf{X} the objective function is not strongly-convex.
4. Consider $n = 100$ and $d = 40$. Generate the data matrix $\mathbf{X} \in \mathbb{R}^{n \times d}$ and the response $\mathbf{y} \in \mathbb{R}^n$, for example, using the python API `numpy.random.randn`. Then calculate the exact solution

$\mathbf{w}^* = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}$ of (1). Use python API to calculate the minimum and maximum eigenvalue of the Hessian matrix, and derive the upper bound on the learning rate η in gradient descent. Let us denote this theoretical bound by η_0 . Run GD on the data set with 6 choices of learning rate: $\eta \in \{0.01\eta_0, 0.1\eta_0, \eta_0, 2\eta_0, 20\eta_0, 100\eta_0\}$. Plot the curve of “ $\|\mathbf{w}^t - \mathbf{w}^*\|_2$ v.s. t ” for $1 \leq t \leq 100$ and summarize your observation. Note that you can start GD with $\mathbf{w}^0 = \mathbf{0}$.

5. Consider $n = 100$ and $d = 200$, and generate \mathbf{X} and \mathbf{y} . What happens when you are trying to calculate the closed-form solution \mathbf{w}^* ? In this case, can we still apply GD? If yes, derive the theoretical bound η_0 and run GD with 6 different η as before. Plot the curve of “ $F(\mathbf{w}^t)$ v.s. t ” for $1 \leq t \leq 100$ and summarize your observation.