

3 Gradient Calculation

Suppose \mathbf{x} and y are known, $\mathbf{w} \in \mathbb{R}^d$ is a column vector. Consider the following functions that have been broadly used in machine learning:

- Sigmoid function $F(\mathbf{w}) = 1/(1 + e^{-\mathbf{x} \cdot \mathbf{w}})$;
- Hinge loss $F(\mathbf{w}) = \max\{1 - y\mathbf{x} \cdot \mathbf{w}, 0\}$;
- ℓ_1 -norm $F(\mathbf{w}) = \|\mathbf{w}\|_1$.

1. Use python to plot their curves for the case $d = 1$. You can set $x = y = 1$;
2. Derive their gradient or subgradients for a general $d > 0$.

3.1 Implementation

Gradient descent is typically used to solve a general optimization problem

$$\min_{\mathbf{w} \in \mathbb{R}^d} F(\mathbf{w}).$$

It starts from an arbitrary point \mathbf{w}^0 and gradually refines the solution as

$$\mathbf{w}^t \leftarrow \mathbf{w}^{t-1} - \eta \cdot \nabla F(\mathbf{w}^{t-1}).$$

Fix $d = 1$, i.e., the variable \mathbf{w} is a scalar. Further fix $w^0 = 1$.

1. Consider $F(w) = \frac{1}{2}w^2$. For each learning rate

$$\eta \in \{10^{-4}, 10^{-3}, 0.01, 0.1, 0.5, 1, 2, 5, 10, 100\},$$

calculate the sequence $\{w^t\}_{t=1}^{1000}$ generated by GD and plot the curve “ $|w^t|$ v.s. t ”.

2. Consider $F(w) = \frac{1}{2}w^4$. For each learning rate

$$\eta \in \{10^{-4}, 10^{-3}, 0.01, 0.1, 0.5, 1, 2, 5, 10, 100\},$$

calculate the sequence $\{w^t\}_{t=1}^{1000}$ generated by GD and plot the curve “ $|w^t|$ v.s. t ”.

4 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:

$$\min_{\mathbf{w} \in \mathbb{R}^d} F(\mathbf{w}) := \frac{1}{2} \|\mathbf{y} - \mathbf{X}\mathbf{w}\|^2, \quad (4.1)$$

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 (6.4) is a convex program.
2. Note that (6.4) is equivalent to the following:

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

in the sense that any minimizer of (6.4) is also an optimum of the above, and vice versa. State why we stick with the least-squares formulation.

3. State when the objective function is strongly-convex and when it is not.
4. Fix $n = 1000$ and increase d from 20 to 500, with a step size 20. For each problem size (n, d) , 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 (6.4) and record the computation time. Plot the curve of “time v.s. d ” and summarize your observation.
5. Consider $n = 100$ and $d = 40$. Again, generate \mathbf{X} and \mathbf{y} , and calculate the optimal solution. 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 (see the slides for the bound). 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 “ $\|\mathbf{w}^t - \mathbf{w}^*\|$ v.s. t ” for $1 \leq t \leq 100$ and summarize your observation. Note that you can start GD with $\mathbf{w}^0 = \mathbf{0}$.
6. 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.