

Notes on Optimization

DSBA Mathematics Refresher 2024

Abstract

Introduction to Optimization

Optimization involves finding the best solution from a set of possible solutions. Specifically, it entails finding the value of x that minimizes or maximizes a function $f(x)$.

Maximization and Minimization

Maximization of $f(x)$ can be transformed into a minimization problem by considering $-f(x)$.

$$\text{maximize } f(x) \equiv \text{minimize } -f(x)$$

Continuous vs. Discrete Optimization

In this session, we will focus on continuous optimization, as it is predominantly used in data science. For your mathematical culture, it's good to know that both exist. Optimization algorithms on continuous cases are quite different to the ones used in discrete cases.

Continuous Optimization: The variables can take any value within a given range, typically \mathbb{R} . *Example:* Linear regression, where parameters can take any real value.

Discrete Optimization: The variables can only take on discrete values, typically \mathbb{Z} or \mathbb{N} . *Example:* Integer programming, where solutions are restricted to integers.

Optimization Algorithms

What algorithms can you think of for solving optimization problems?

[Interactive session]

Here, are methods presented in class:

Grid Search

Grid Search is a brute-force method that evaluates the function at a grid of points covering the domain. It is simple but computationally expensive, especially in high dimensions. If your domain is unbounded (such as \mathbb{R}), you will need to bound your grid search.

Dichotomy (Bisection Method)

The Bisection Method is fitted for one-dimensional optimization (although it is possible to extend it to multi-dimension). It repeatedly bisects an interval and selects the sub-interval in which the function changes sign. It is effective for finding roots but can be extended to optimization.

Gradient Descent

Gradient Descent is an iterative method used for finding local minima of a function. It updates the parameters in the opposite direction of the gradient of the function at the current point.

The update rule is: $x_{k+1} = x_k - \alpha \nabla f(x_k)$, where α is the learning rate and $\nabla f(x_k)$ is the gradient of f at x_k .

In one dimension, $\nabla f(x_k)$ is simply the derivative of f .

Algorithms in Practice

Grid Search

- **Step 1:** Define the grid over the domain.
- **Step 2:** Evaluate the function at each grid point.
- **Step 3:** Select the point with the best function value.

Advantages: Simple and straightforward.

Disadvantages: Computationally expensive, especially in high dimensions. It also needs a bounded domain.

Bisection Method

- **Step 1:** Choose initial interval $[a, b]$ to explore.
- **Step 2:** Compute the midpoint $c = \frac{a+b}{2}$.
- **Step 3:** Determine the "best"¹ sub-interval $[a, c]$ or $[c, b]$.
- **Step 4:** Repeat until the interval is sufficiently small.

¹"best" means the one most likely to contain the minimum.

Advantages: Fun to implement.

Disadvantages: Extending the technique to multi dimensional problem is complicated. It also needs a bounded domain.

Gradient Descent

- **Step 1:** Initialize x_0 .
- **Step 2:** Compute the gradient $\nabla f(x_k)$.
- **Step 3:** Update $x_{k+1} = x_k - \alpha \nabla f(x_k)$.
- **Step 4:** Repeat until convergence (i.e., $\|\nabla f(x_k)\|$ is small).

Advantages: Efficient for high-dimensional problems, widely used in machine learning.

Disadvantages: May converge to local minima, requires tuning of the learning rate.

Conclusion

Continuous optimization is a crucial tool in data science. It's a good exercise to test various optimization techniques, and you are encourage to try other ones on you own (genetic algorithms are interesting as well, although only used in niche cases). In practice, gradient descent (and variations of it, such as Adam) are used in the vast majority of the cases.

Linear Regression

This section will look at 1D-1D linear regression so that formulas maintain a manageable weight. The same principle can be applied to ND -1D linear regressions, and repeating the process for each output dimension, it is possible to extend it to ND - MD cases.

Given a set of data points $(x_k, y_k) \in \mathbb{R}^2$ for $1 \leq k \leq N$, we want to find the line $y = ax + b$ that best fits the data. The function $y = f(x)$ is a linear model with parameters a and b , where a represents the slope and b represents the intercept.

Least Squares Method

To define the "best fit" for a line, we use the sum of errors squared. The idea is to minimize the sum of the squared differences between the observed values y_k and the predicted values $f(x_k)$. This sum will be our cost function (or loss):

$$\mathcal{L}(a, b) = \sum_{k=1}^N (y_k - (ax_k + b))^2$$

Our objective is to minimize $\mathcal{L}(a, b)$ with respect to a and b .

Derivation of a and b

To minimize $\mathcal{L}(a, b)$, we calculate the partial derivatives of \mathcal{L} with respect to a and b .

First, the partial derivative of \mathcal{L} with respect to a is:

$$\begin{aligned}\frac{\partial \mathcal{L}(a, b)}{\partial a} &= \frac{\partial}{\partial a} \sum_{k=1}^N (y_k - (ax_k + b))^2 \\ &= \sum_{k=1}^N \frac{\partial}{\partial a} (y_k - ax_k - b)^2 \\ &= \sum_{k=1}^N 2(y_k - ax_k - b)(-x_k) \\ &= -2 \sum_{k=1}^N x_k (y_k - ax_k - b)\end{aligned}$$

$\mathcal{L}(a, b)$ is minimum if $\frac{\partial \mathcal{L}(a, b)}{\partial a} = 0$, so we want:

$$\sum_{k=1}^N x_k y_k - a \sum_{k=1}^N x_k^2 - b \sum_{k=1}^N x_k = 0$$

Next, the partial derivative of $\mathcal{L}(a, b)$ with respect to b is:

$$\begin{aligned}\frac{\partial \mathcal{L}(a, b)}{\partial b} &= \frac{\partial}{\partial b} \sum_{k=1}^N (y_k - (ax_k + b))^2 \\ &= \sum_{k=1}^N \frac{\partial}{\partial b} (y_k - ax_k - b)^2 \\ &= \sum_{k=1}^N 2(y_k - ax_k - b)(-1) \\ &= -2 \sum_{k=1}^N (y_k - ax_k - b)\end{aligned}$$

Setting $\frac{\partial \mathcal{L}(a, b)}{\partial b} = 0$:

$$\sum_{k=1}^N y_k - a \sum_{k=1}^N x_k - Nb = 0$$

Solving the System of Equations

We now have a system of two linear equations to solve for a and b :

$$\begin{aligned}\sum_{k=1}^N x_k y_k &= a \sum_{k=1}^N x_k^2 + b \sum_{k=1}^N x_k \\ \sum_{k=1}^N y_k &= a \sum_{k=1}^N x_k + Nb\end{aligned}$$

[Details left to the reader]

Setting $\bar{x} = \frac{1}{N} \sum_{k=1}^N x_k$ and $\bar{y} = \frac{1}{N} \sum_{k=1}^N y_k$ (i.e. the means of the x and y values, respectively):

$$\begin{aligned}a &= \frac{\left(\sum_{k=1}^N x_k y_k\right) - N\bar{x}\bar{y}}{\sum_{k=1}^N x_k(x_k - \bar{x})} \\ b &= \bar{y} - a\bar{x}\end{aligned}$$