

Prof. Sören Petrat, Constructor University

Lecture notes from Spring 2025

6. Multivariable Calculus6.2 Optimization in Many Variables

Topic for Week 12B: Least Squares and Gradient Descent

Let us make another connection to the Elements of Linear Algebra class from last semester.

There, in Week 12B, we considered the Least Square problem:

Given an $m \times n$ matrix A with $m > n$ and a vector $b \in \mathbb{R}^m$, find $x \in \mathbb{R}^n$ that minimizes $\|Ax - b\|$.

In linear Algebra, we used the QR decomposition to solve this.

Now, let us use Calculus to solve this problem.

Let us define $f(x) = \|Ax - b\|^2$ which has the same minimum as $\|Ax - b\|$.

$$\begin{aligned} f(x) &= (Ax - b)^T (Ax - b) = (x^T A^T - b^T)(Ax - b) \\ &\quad \text{recall: } \|x\|^2 = x^T x \\ &\quad (A^T = \text{transpose of } A) \\ &= x^T A^T A x - x^T A^T b - b^T A x + b^T b \\ &\quad \text{recall } (AB)^T = B^T A^T \\ &= x^T (b^T A)^T \\ &= x^T A^T b \end{aligned}$$

$$\begin{aligned} &= x^T A^T A x - 2x^T A^T b + b^T b \\ &= \sum_{j \in R} x_j (A^T A)_{jj} x_j - 2 \sum_j x_j (A^T b)_j + b^T b \end{aligned}$$

Minimization: $0 = (\nabla f)(x)$ i.e., $\frac{\partial f}{\partial x_i} = 0 \quad \forall i = 1, \dots, n$.

$$\Rightarrow \frac{\partial f}{\partial x_i} = \sum_e (A^T A)_{ie} x_e + \underbrace{\sum_j x_j (A^T A)_{ji}}_{= (A^T A)_{ij}^T x_j} - 2 (A^T b)_i$$

$$\Rightarrow 0 = \nabla f = 2A^T A x - 2A^T b \Rightarrow A^T A x = A^T b$$

\Rightarrow The minimum occurs when $A^T A x = A^T b$ (which is an $n \times n$ system of linear equations that still needs to be solved)

Note: Here our function was very simple, hence we could do the minimization exactly.

But in many applications, especially in Machine learning, we cannot find minima exactly. Hence we need numerical algorithms.

Let us briefly discuss the important method of gradient descent. The goal is to numerically find minima of $f: \mathbb{R}^n \rightarrow \mathbb{R}$.

We start from our observation that ∇f points in the direction of the strongest increase. Hence we could use the following

Algorithm:

- Choose an initial point $a_0 \in \mathbb{R}^n$.
- Choose a step size $\eta > 0$ and go a step η in the direction of $-\nabla f$ (steepest descent):

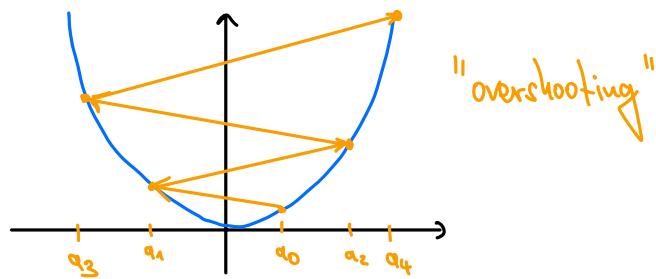
$$a_1 = a_0 - \eta (\nabla f)(a_0)$$
- Repeat:
$$a_{n+1} = a_n - \eta (\nabla f)(a_n)$$
- Stop when $\|(\nabla f)(a_n)\| < \varepsilon$ for some tolerance ε .

Note: In the context of Machine learning, the step size η is called "learning rate".

Applications and implementations of this method are discussed in other classes. Here, let us discuss some potential problems and generalizations.

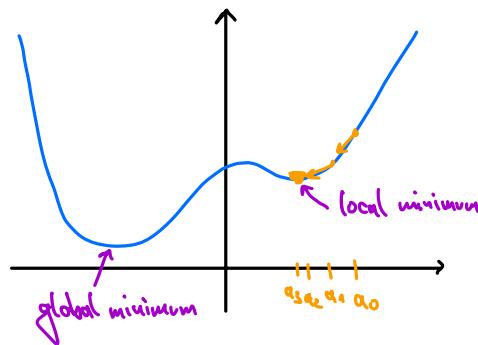
Potential problems:

- Learning rate:
 - If η is chosen too large, we might not find the minimum, or the algorithm might even diverge.

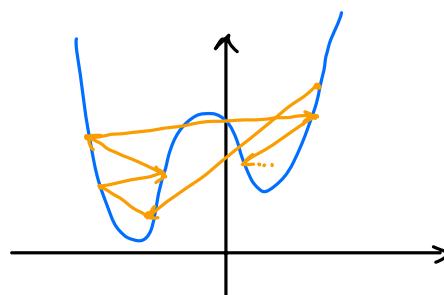


- If η is chosen too small, the convergence might be very slow.

- Local vs. global minima: The algorithm might get stuck in a local minimum instead of a global minimum, especially if the learning rate is chosen too small.



- Related unlucky situations:



Generalizations:

- Making the learning rate n-dependent: For example, decreasing the learning rate with increasing number of iterations mitigates the risk of overshooting the minimum. E.g., $\eta_n = \frac{1}{1+n}$
- Stochastic gradient descent: Often, $f(x) = \frac{1}{N} \sum_{i=1}^N f_i(x)$ with N large, so evaluating ∇f is computationally expensive. Here, instead of using the iteration $a_{n+1} = a_n - \eta \frac{1}{N} \sum_{i=1}^N (\nabla f_i)(a_n)$, we choose one i randomly in each step and perform the iteration $a_{n+1} = a_n - \eta \underbrace{(\nabla f_i)}_{\text{line}}(a_n)$.

Advantages:

- computationally cheap
- might be able to escape local minima

Or we could take an average over a small subset $M \subset \{1, \dots, N\}$: $\frac{1}{|M|} \sum_{j \in M} (\nabla f_j)(a_n)$

Disadvantages:

- might need more steps to reach minimum
- might "jump around" near an actual minimum