

What are the limitations of derivative-based models for optimization in machine learning?

Faris Chaudhry

August 5, 2021

Abstract

Most machine learning problems can be transposed into optimization problems with the goal being finding the global minima or maxima to minimize loss or maximize potential. Each of the main learning methodologies (namely supervised, unsupervised and reinforcement) along with the models to represent and solve optimization problems have limitations - computationally and theoretically - that have to be identified and mitigated against to create an effective model. The focus here are objective functions that are continuously differentiable and thus derivative-based solutions are used.

Contents

1	Introduction to Machine Learning and Optimization	1
1.1	What is Machine Learning?	1
1.2	Prerequisite Conditions for Derivative-Based Optimization	2
1.2.1	Continuity and Differentiability	2
1.2.2	Concavity and Convexity	3
2	Supervised Learning	4
2.1	Application of Supervised Learning	4
2.1.1	Classification	4
2.1.2	Regression	4
2.2	General Optimization of Supervised Learning Problems	5
2.3	Limitations of Supervised Learning	5
3	Unsupervised Learning	6
3.1	Application of Unsupervised Learning	6
3.1.1	Clustering	6
3.1.2	Density Estimation	7
3.2	Limitations of Unsupervised Learning	8
4	Reinforcement Learning	9
4.1	Optimization in Reinforcement Learning	9
4.1.1	Limitations of Reinforcement Learning	10
5	Mathematical Models for Optimization	11
5.1	Constraints on the Objective Function	11
5.2	Methods for Finding Extrema	11
5.2.1	Jacobian	11
5.2.2	Hessian	12
5.2.3	Higher-Order Derivative Tests	12
5.2.4	Iterative Methods	13

Chapter 1

Introduction to Machine Learning and Optimization

1.1 What is Machine Learning?

Machine learning (ML) is a subfield of artificial intelligence (AI) which, broadly speaking, is the use of computational methods and models to improve performance and predictions through experience [2, p. 1]. Unlike humans, this learning is based entirely on data and statistics and experience is gained through interaction with a training set of data or an environment of some kind.

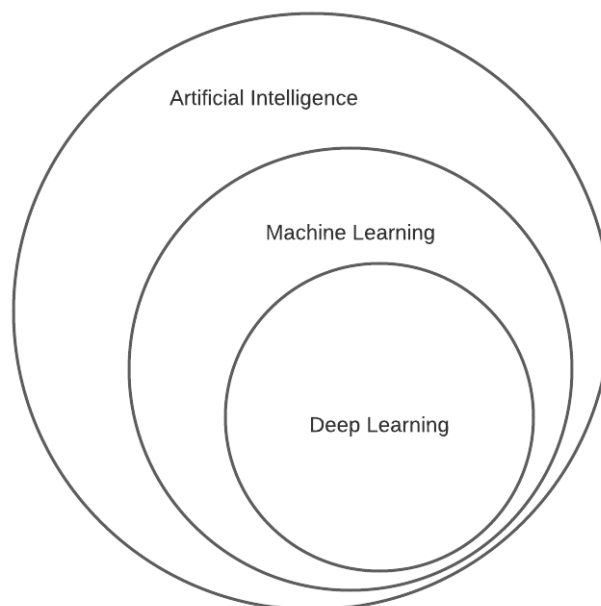


Figure 1.1: subfields of AI

There are 3 primary categories (supervised, unsupervised and reinforcement) of learning philosophies for ML models, with other hybrid models being combinations of these. Each type of learning lends to itself to certain types of problems due to the limitations that each one has. Supervised learning is used for classifying images and extrapolating data. Unsupervised learning takes raw data and finds patterns such as the overall distribution or groups with similar attributes. Reinforcement is used in complex systems which many changing variables that would be computationally difficult to solve otherwise, like chess.

1.2 Prerequisite Conditions for Derivative-Based Optimization

Optimization revolves around minimizing the loss or maximizing the value of a function. In the context of ML, the process of optimizing is vital to ensure modelling produces the greatest accuracy. The goal is to optimize an objective function, which is the representation of the variables being simulated. The solution to the objective function will be either a minimum (minima) or maximum (maxima) point (collectively called the set of extrema) as this is when the value of a function is highest or lowest.

The derivative is a linear approximation (tangent) to a function at a point. Suppose there was a function $f(x)$ then, intuitively, the derivative with respect to x would be the how much the value of $f(x)$ changed with a small nudge in the x direction. It is important to note that, by Fermat's theorem on stationary points [12], all critical points (extrema and saddle points) have first derivative equal to 0. Visually, this is because the tangent to any turning point will have a gradient of 0. See the function $y = x^2$ (fig. 1.2) which has a minima at $(0,0)$.

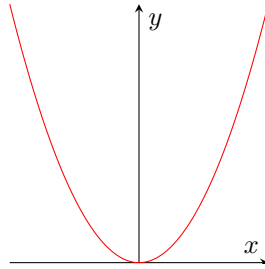


Figure 1.2: graph of x^2

So the objective becomes to find all critical point locations and the nature of these points - what kind of critical point it is.

1.2.1 Continuity and Differentiability

The most essential requirement to using derivative-based methods will be that the objective function must be continuous and twice-differentiable (the derivative of the function must also be differentiable) over the interval that contains the solution.

This is because to find the location and nature of critical points, the first and second derivative of a function are required [11].

For a function $f(x)$ to be continuous over the interval $I = [a, b]$

$$\forall k \in I, \lim_{x \rightarrow k} f(x) = f(k) \quad (1)$$

This means that, given any number in the interval, as x approaches that number it would be equal to putting the number into the function. This prevents any discontinuity since the limit wouldn't exist at discontinuous points. In fig. 1.3, $\lim_{x \rightarrow 0+} = +\infty$ and $\lim_{x \rightarrow 0-} = -\infty$. These values contradict meaning the limit isn't defined.

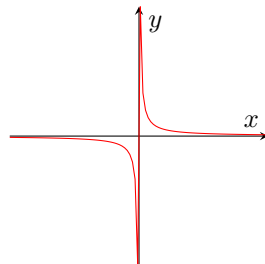


Figure 1.3: graph of $\frac{1}{x}$

For a single-valued function, $f(x)$, the derivative, $f'(x)$ exists iff the following limits exists.

$$\frac{df}{dx} = \lim_{\Delta x \rightarrow 0} \frac{f(x + \Delta x) - f(x)}{\Delta x} \quad (2)$$

However, most objective functions will be multi-valued to account for all the variables so this definition must be extended. This is the same principle but a nudge in a specific direction. Suppose there is a function $f(x_1, \dots, x_i)$ then the derivative with respect to a certain variable, x_n , will be.

$$\frac{\partial f}{\partial x_n} = \lim_{x \rightarrow 0} \frac{f(x_1, \dots, x_n + \Delta x, \dots, x_i) - f(x_1, \dots, x_n, \dots, x_i)}{\Delta x} \quad (3)$$

In practice these rigorous definitions are not used but the concept of continuous differentiability is important.

- The first and second derivative must exist for an objective function to be solvable in this method, which is the major limiting factor. Although derivative-free methods do exist, they tend to be approximations of the exact values and heuristic in theory.
- Although many functions discontinuities, like asymptotes or singularities, many times these are removable either by defining an interval without them or assigning an arbitrary value at a point for continuity.

1.2.2 Concavity and Convexity

When a function has only 1 minima or maxima over an interval it becomes much easier to find the global minimum or maximum due the lack of a need to check which point is a local extremum and which is the global extremum. Functions like these are called convex and concave where convex functions have a minimum point and concave functions have a maximum point. A convex function [7] can visually be described as having all its points below a line segment drawn between (fig. 1.4) any 2 points while a concave function has all points above.

It is important to note that concavity and convexity are not opposites. A function can be concave, non-concave, convex or non-convex. In addition, reflecting a function in the x -axis will reverse its concavity or convexity. Suppose $f(x)$ is convex then $-f(x)$ is concave.

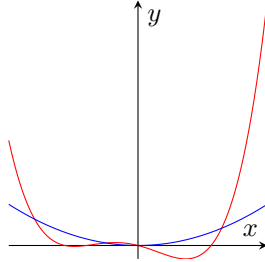


Figure 1.4: convex (blue) and non-convex (red)

If the objective function is a non-convex or non-concave function this doesn't prevent the use of derivative-based optimization. However it does restrict the range of methods that can be used to find the global solution. For example, iterative methods to find extrema might not always work since they could get closer to the local extrema while neglecting other possible values. Moreover, it will increase the complexity of the problem computationally since there will be range of possible global extrema that have to be checked - which can be particularly difficult when certain derivative tests are inconclusive.

Chapter 2

Supervised Learning

2.1 Application of Supervised Learning

The philosophy of supervised learning is to use labelled training data to map between an input vector and a target vector. In this case, the model is given data with input variables and the correct associated target values corresponding with them [1, p. 105]. Effectively, the model is creating a pattern out of which inputs cause certain outputs so that, given new inputs, the correct outputs can be predicted.

Supervised learning problems are split into 2 distinct categories: classification problems and regression problems.

2.1.1 Classification

Classification problems are about predicting the class labels of an object. A common example of classification is assigning an digit label to a handwritten digit. However, these objects could be anything that can be labelled such as sentences or sounds.

Let x_n be a feature/parameter of the object and l_n be a label where $n \in \mathbb{Z}$
Then a general classification function can be described as the mapping:

$$[x_1, x_2, \dots] \mapsto [l_1, l_2, \dots]$$

Given a particular feature vector (a vector of the parameters of the classification function) the goal is to assign a set of class labels.

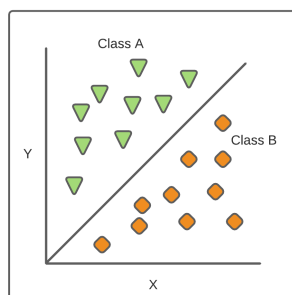


Figure 2.1: example of classification

2.1.2 Regression

Regression problems involve predicting a numerical value from the feature vector of an object. For example, given many variables about a stock (past history), predict the future value of the stock.

Let x_n be a feature/parameter of the object and k be the numerical value associated with it where $n \in \mathbb{Z}$. Then a general regression function can be described as the mapping:

$$[x_1, x_2, \dots] \mapsto k$$

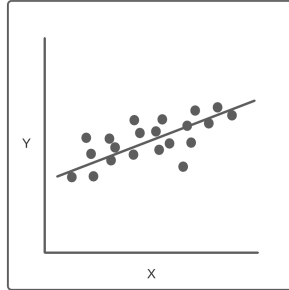


Figure 2.2: example of regression

2.2 General Optimization of Supervised Learning Problems

The optimization of a supervised learning problem is to minimize the average of the loss function using the training samples. This produces the most accurate approximation to the underlying function to extrapolate values.

The general equation [6, p. 3] for this can be written as:

$$\min_{\theta} \frac{1}{N} \sum_{i=1}^N L(y^i, f(x^i, \theta)) \quad (4)$$

where N is the number of training samples, θ is the parameter of the mapping function, x^i is a feature vector and y^i is the array of labels associated with that feature vector.

The problem with using training samples is that the resulting function might be over fitted to the given data. This would mean that, although the model is accurate for the training data it has been given, accuracy is reduced on new objects. The method to deal with this is through a regularization item, λ :

$$\min_{\theta} \frac{1}{N} \sum_{i=1}^N L(y^i, f(x^i, \theta)) + \lambda \|\theta\|_2^2 \quad (5)$$

Regularization fundamentally discourages learning complex models and will be covered more in depth when talking about limitations.

2.3 Limitations of Supervised Learning

Models don't generalize well from observed, training data to unseen data [14]. The accuracy might be near perfect on training data while being poor on unseen data somewhat mimicking the model memorizing the training data without grasping the mapping function.

Chapter 3

Unsupervised Learning

3.1 Application of Unsupervised Learning

Unsupervised learning is different to supervised learning in the way that it uses unlabelled data [1, p. 105]; instead of learning from a mapping of inputs to a known output, the model is given only the inputs to learn from and has to make sense of the data without guidance. As a result of this, unsupervised learning revolves around extracting relationships from the data without the inherent human biases caused by choosing the correct output beforehand.

Unsupervised learning problems strive to solve 1 of 2 problems: finding clusters of similar data and summarizing the distribution of the data (density estimation).

3.1.1 Clustering

Unlike classification, where the classes are predefined, clustering requires the model to define its own cluster of data based on the similarities of the features.

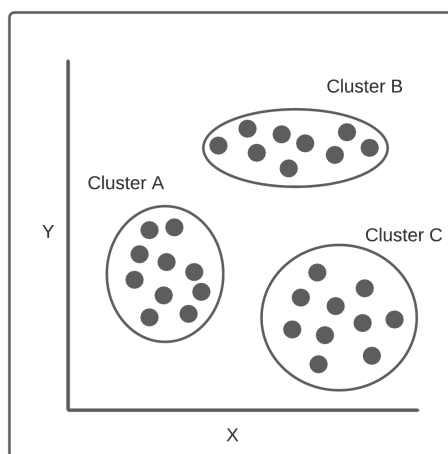


Figure 3.1: example of clustering

Optimization will involve making the variance of each cluster as small as possible; variance will be equivalent to the distance from the center of the cluster. This can be done by iterating through each cluster and checking the distance from the centre of the cluster for each sample.

$$\min_s \sum_{k=1}^K \sum_{x \in S_k} \|x - \mu_k\|_2^2 \quad (6)$$

where s is the variance, K is the number of clusters, S_k is the set of samples for that cluster, μ_k is the centre of a cluster [6, p. 3-4].

3.1.2 Density Estimation

The assumption is that there exists some probability distribution to describe the relationship between the variables [4]. Density estimation is a useful asset in modelling to estimate the properties of a given data set (variance, skewness, type of distribution).

Suppose there exists a set of continuous random variables, (x_1, \dots, x_n) , then there is a probability distribution that the set models, $P(x_1, \dots, x_n)$. The goal is to find a continuous probability density function (PDF) that can describe the mapping: $\{x_1, \dots, x_n\} \rightarrow P(x_1, \dots, x_n)$. The assumption that this will PDF will be continuous is valid since our objective function has the precondition of being continuous.

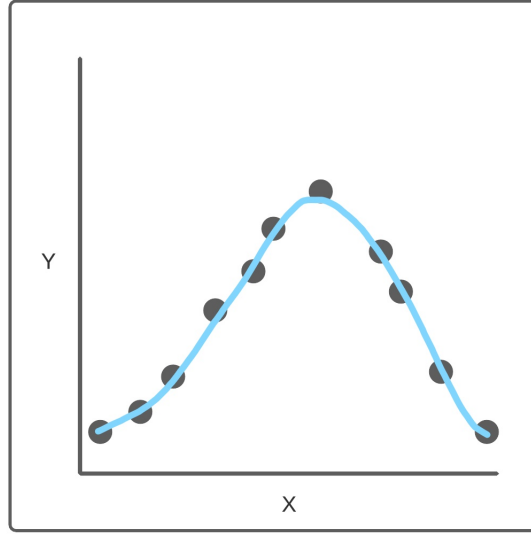


Figure 3.2: example of density estimation

Getting an exact function approximation is not only unlikely but also suboptimal. As more data is added the approximation will get more complex and usually closer to the actual value however too much data increases the risk of overfitting, especially if there is abnormal data. Instead, the focus is to maximize the likelihood that a predicted PDF is correct.

In statistics, the likelihood function measure the goodness of fit between the data and the proposed PDF so maximizing it finds the PDF which has the highest probability of fitting the data. Principally, a likelihood function is written as $p(x^i, \theta)$ where x^i is a datapoint and θ is a PDF. However, many data sets have large numeric ranges that can limit the effectiveness of this function so it is easier to work with the logarithmic likelihood function, $\ln p(x^i, \theta)$ or $\ell(x^i, \theta)$.

It is important to note that

$$\max p(x^i, \theta) = \max \ell(x^i, \theta) = \widehat{\ell}(x^i, \widehat{\theta}) \quad (7)$$

Since \ln is strictly increasing then maximizing the logarithmic likelihood function also maximizes the likelihood function. Thus general density estimation [6, p. 4] is:

$$\max \sum_{i=1}^N \ell(x^i, \theta) \quad (8)$$

where N is the number of training samples and x^i is a particular feature vector of a sample.

3.2 Limitations of Unsupervised Learning

Chapter 4

Reinforcement Learning

Reinforcement learning entails an agent interacting with an environment rather than using a traditional dataset [1, p. 105]. The learning process is done through trial and error until there is a feedback loop between the environment and the agent's experience. Each situation or state that the environment can have should be mapped to an action that the agent can take to maximize reward.

This game theory approach assigns a payoff yield to each action that encourages or discourages certain actions. A byproduct of this is to consider how a length of time can affect the payoff of an action. For example, maybe an action has a high reward when considering only the next state but when looking at the next 10 states its value is lower. This can be seen in chess when, as the depth of the AI get higher, certain moves become worse since they compromise future positions.

The dynamic nature of reinforcement learning is useful in complex systems where computationally working out the whole game would be inefficient if not impossible. Consider chess [3] where the number of possible games is estimated as 10^{123} (Shannon's number) and is too big to compute after even 5 turns.

Number of half-moves	Number of Possible Games
1	20
2	400
3	8,902
4	197,281
5	4,865,609

Figure 4.1: the exponential growth of possible chess games

4.1 Optimization in Reinforcement Learning

A policy function, $\pi(s)$ maps a state, s , to an action, a , to select the best course of action from the set of all actions, A , that the agent can take in each situation from the set of possible situations, S .

$$\pi(s) : s \rightarrow a \text{ where } s \in S, a \in A$$

By maximizing the expected value of this function [6, p. 4], the agent will select the best action to be performed in each state, maximizing the payoff of the actions.

$$\max_{\pi(s)} = \mathbf{E} \left[\sum_{k=1}^T \gamma^k r_{t+k} | S_t = s \right] \quad (9)$$

where T is the time horizon, γ is the discount factor, r is the reward function with respect to the turn and the time into the future being considered, S_t is a given state.

If the game was infinite - or at least has no predetermined stopping point - then a simple solution is to work out the $\lim_{T \rightarrow \infty}$ and truncate the series at some point to get a good approximation. Using this method of truncation would allow the depth to be changed (how far into the future turns are considered).

The discount factor, γ , is an value that prioritizes instant reward over a future reward [5]. If there is a constant risk which may cause failure to realize the reward then that future reward should have a decreased payoff to compensate for the risk. For example, suppose there is a 50% chance (implying $\gamma = 0.5$) that the game ends after every turn and you could choose either a payoff of 1 unit after 1 turn or a payoff of 10 units after 5 turns. Then the expected value of option 1 is $0.5^1 * 1$ which is 0.5 and option 2 is $0.5^5 * 10$ which is 0.3125 so option 1 is statistically better.

4.1.1 Limitations of Reinforcement Learning

Chapter 5

Mathematical Models for Optimization

5.1 Constraints on the Objective Function

Constraints are conditions that the solutions to the objective function must satisfy. In many cases, there are certain restrictions that should be added due to computational and resource limitations.

There are 3 types of constraints that must be considered:

- inequality constraints such as $x \geq k$.
- equality constraints such as $x = k$
- data type constraints such as x is an integer ($x \in \mathbb{Z}$)

Data type constraints are normally easy to deal with; changing with values the model checks or adding a conditional statement can be enough in many cases.

Equality constraints can be appended onto the objective function using the Lagrange multiplier [10]. Let $f(x)$ be an objective function and $c_n(x) = k_n$ be an equality constraint. Then the general lagrange multiplier for n constraints would be:

$$\mathcal{L}(x) = f(x) - \lambda_1(c_1(x) - k_1) - \cdots - \lambda_n(c_n(x) - k_n) \quad (10)$$

The solution to the original objective function, $f(x)$ will be a saddle point on the lagrange multiplier. One problem with this approach is that, firstly, for each constraint added there will be an extra partial derivative that needs to be computed and thus a more complex simultaneous equation. In addition, inequality constraints ($c(x) \geq k$) are not supported through this method. However, there is a generalization to the lagrange multiplier called the Karush–Kuhn–Tucker conditions (KKT).

5.2 Methods for Finding Extrema

5.2.1 Jacobian

Previously, we ascertained that the locations of the critical points required the first derivative of the function to be 0 ($f' = 0$). For a single-valued function this is equivalent to solving

$$\frac{dy}{dx} = 0$$

However for multi-valued functions a critical point will require all partial derivatives to be 0 at that point. The Jacobian [13] is a vector of the partial derivatives of the function that gives a vector of the

gradients at that point.

$$J_f = \nabla f = [\partial_{x_1} f, \dots, \partial_{x_n} f] \quad (11)$$

This can be solved by either setting each partial derivative to 0 and solving it simultaneously or by finding the values when the value of the jacobian is 0.

$$\partial_{x_1} f = \dots = \partial_{x_n} f = 0 \text{ or } \|J\| = 0$$

Sometimes the resulting equations might be difficult to solve (polynomials of order greater than 4 for example) and brute force or numerical methods might need to be used.

5.2.2 Hessian

After finding the locations of the critical points, higher order derivative tests must be used to determine the nature of the points. The Hessian [9] is like the 'jacobian of the jacobian' and is a representation of the change in the gradient.

$$H_f = \nabla(\nabla f) = \begin{bmatrix} \frac{\partial^2 f}{\partial x_1^2} & \frac{\partial^2 f}{\partial x_1 \partial x_2} & \dots & \frac{\partial^2 f}{\partial x_1 \partial x_n} \\ \frac{\partial^2 f}{\partial x_2 \partial x_1} & \frac{\partial^2 f}{\partial x_2^2} & \dots & \frac{\partial^2 f}{\partial x_2 \partial x_n} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial^2 f}{\partial x_n \partial x_1} & \frac{\partial^2 f}{\partial x_n \partial x_2} & \dots & \frac{\partial^2 f}{\partial x_n^2} \end{bmatrix} \quad (12)$$

Due to the commutative nature of second partial derivatives, computation time is significantly cut down as only the upper or lower triangular matrix has to be calculated and then mirrored.

$$\frac{\partial^2 f}{\partial_y \partial_x} = \frac{\partial^2 f}{\partial_x \partial_y} \quad (13)$$

5.2.3 Higher-Order Derivative Tests

Using the determinant and trace of the Hessian matrix where λ is an eigenvalue

$$trH = \sum_{i=1}^n H_{ii} = \prod_{\forall i} \lambda_i \text{ and } detH = \sum_{\forall i} \lambda_i \quad (14)$$

the second derivative test is as follows. Let (x_1, x_2, \dots, x_n) be a critical point substituted into the Hessian. Then [11]

- $detH > 0$ and $trH > 0 \implies$ local minimum.
- $detH > 0$ and $trH < 0 \implies$ local maximum.
- $detH < 0 \implies$ saddle point.
- $detH = 0 \implies$ inconclusive test.

Doing this for all critical points will determine which are extrema and comparing between the extrema will find the global minimum or maximum, which will be the solution to the optimization problem. In the case when the test is inconclusive, higher-order tests [8] can be used (although depending on the variables many derivatives would have to be calculated) or inspection.

For complex functions, visual inspection is usually impossible (since displaying something with greater than 3 dimensions is unintuitive) but numerical inspection can be useful, like adding dx_n for each variable and comparing the values. For example, if the neighbourhood of a point has values greater than the point then it must be a minima.

The Hessian matrix is not without limitations; the computing time required to work out the Hessian scales greatly as each new variable is added and, for each critical point, the eigenvalues of the matrix will have to be worked out to find the trace and determinant. To mitigate this there are algorithms which can approximate the Hessian with the only condition that the matrix is invertible (the determinant is not 0).

5.2.4 Iterative Methods

When exact solutions are too difficult to calculate either due to too many variables or the Hessian matrix being too large to reasonably store, iterative methods can be used as numerical approximations. Since the approximation will converge to the answer as the number of iterations approaches infinity, many iterative methods are far superior, in practice, because of the lower resource cost relative to the accuracy of the approximation.

List of Figures

1.1	subfields of AI	1
1.2	graph of x^2	2
1.3	graph of $\frac{1}{x}$	2
1.4	convex (blue) and non-convex (red)	3
2.1	example of classification	4
2.2	example of regression	5
3.1	example of clustering	6
3.2	example of density estimation	7
4.1	the exponential growth of possible chess games	9

References

- [1] Yoshua Bengio, Ian Goodfellow, and Aaron Courville. *Deep learning*. Vol. 1. MIT press Massachusetts, USA: 2017.
- [2] Mehryar Mohri, Afshin Rostamizadeh, and Ameet Talwalkar. *Foundations of machine learning*. MIT press, 2018.
- [3] Claude E Shannon. “XXII. Programming a computer for playing chess”. In: *The London, Edinburgh, and Dublin Philosophical Magazine and Journal of Science* 41.314 (1950), pp. 256–275.
- [4] Simon J Sheather. “Density estimation”. In: *Statistical science* (2004), pp. 588–597.
- [5] Peter D Sozou. “On hyperbolic discounting and uncertain hazard rates”. In: *Proceedings of the Royal Society of London. Series B: Biological Sciences* 265.1409 (1998), pp. 2015–2020.
- [6] Shiliang Sun et al. “A survey of optimization methods from a machine learning perspective”. In: *IEEE transactions on cybernetics* 50.8 (2019), pp. 3668–3681.
- [7] Lieven Vandenbergh and Stephen Boyd. *Convex optimization*. Vol. 1. Cambridge University Press Cambridge, 2004.
- [8] Eric W. Weisstein. *Extremum Test*. From MathWorld—A Wolfram Web Resource. URL: <https://mathworld.wolfram.com/ExtremumTest.html>.
- [9] Eric W. Weisstein. *Jacobian*. From MathWorld—A Wolfram Web Resource. URL: <https://mathworld.wolfram.com/Hessian.html>.
- [10] Eric W. Weisstein. *Lagrange Multiplier*. From MathWorld—A Wolfram Web Resource. URL: <https://mathworld.wolfram.com/LagrangeMultiplier.html>.
- [11] Eric W. Weisstein. *Second Derivative Test*. From MathWorld—A Wolfram Web Resource. URL: <https://mathworld.wolfram.com/SecondDerivativeTest.html>.
- [12] Eric W. Weisstein. *Stationary Point*. From MathWorld—A Wolfram Web Resource. URL: <https://mathworld.wolfram.com/StationaryPoint.html>.
- [13] Eric W. Weisstein. *Stationary Point*. From MathWorld—A Wolfram Web Resource. URL: <https://mathworld.wolfram.com/Jacobian.html>.
- [14] Xue Ying. “An overview of overfitting and its solutions”. In: *Journal of Physics: Conference Series*. Vol. 1168. 2. IOP Publishing. 2019, p. 022022.