

Machine learning Algorithms

蔡苗¹

2019-04-18

¹Department of Epidemiology and Biostatistics, College for Public Health and Social Justice, Saint Louis University. Email: miao.cai@slu.edu

感谢我的家人的支持。

Acknowledgement

I want to thank my mentor.

目录

Preface	xi
第一章 Introduction	1
第二章 Discrete optimization	3
2.1 Heuristic and metaheuristic methods	3
2.2 Genetic algorithm and simulated Annealing as examples	5
2.3 Constrains	5
第三章 Continuous optimization	7
3.1 First and second derivative	7
3.2 First Order Derivative methods	8
3.2.1 Gradient descent	8
3.2.2 Stochastic gradient descent	9
3.2.3 Coordinate descent	9
3.3 Second Order Derivative methods	10
3.3.1 Iteratively reweighted least squares (IRLS)	10
3.3.2 Newton-Raphson optimization and Fisher Scoring	10
3.3.3 Limited-memory Broyden-Fletcher-Goldfarb-Shanno (L-BFGS)	10
3.3.4 L-BFGS Versus IRLS for GLM's	11
3.4 Close thoughts	11
第四章 Introduction	13
第五章 Introduction	15

第六章 Introduction

17

表格

插图

Preface

This book works as a notebook to summarize the algorithms used in Bayesian inference and machine learning.

第一章 Introduction

第二章 Discrete optimization

Most of the concepts are explain in Chapter 4 Numerical Computation from the Deep Learning book <https://www.deeplearningbook.org/contents/numerical.html>.

The **objective function** allows us to measure how “good” any given solution to the problem is. We seek to maximize or minimize the objective function.

Derivative/gradient based methods keep going “uphill” until they are at the top of the h

2.1 Heuristic and metaheuristic methods

a **heuristic** is a technique designed for solving a problem more quickly when classic methods are too slow, or for finding an approximate solution when classic methods fail to find any exact solution.

Wikipedia

Heuristic methods do not guarantee to find the global optimal solution (best solution)! Instead, they seek to find a **best available solution, given the resource spent looking for it**. A **heuristic method** is geared towards a specific problem.

a **metaheuristic** is a higher-level procedure or heuristic designed to find, generate, or select a heuristic (partial search algorithm) that may

provide a sufficiently good solution to an optimization problem, especially with incomplete or imperfect information or limited computation capacity. Metaheuristics sample a set of solutions which is too large to be completely sampled. Metaheuristics may make few assumptions about the optimization problem being solved, and so they may be usable for a variety of problems

– Wikipedia

A metaheuristic method is like a heuristic, but generalizable to a broad class of problems.

1. Genetic Algorithms (Holland – 1975)
 - Natural selection / genetics based. Popular method.
2. Simulated Annealing (Kirpatrick – 1983)
 - Metallurgy annealing, find lowest energy level!
3. Particle Swarm Optimization (Eberhart Kennedy - 1995)
 - Based on insect behavior, swarming towards optimal location (food). Less common in discrete spaces. originally proposed for continuous spaces.
4. Tabu Search (Al-Sultan – 1999)
 - Search for best neighborhood solution, then find new neighborhood. Prior neighborhoods are forbidden (tabu)

General meta-heuristics traits

- Evaluate many potential optimal solutions.
- Evaluate the fitness of each solution based on a cost (objective) function.
- Use some concept of stochastic (random) movement to generate new solutions from the parameter space.
- Use some set of rules to determine where to move next in the parameter space.
- Declare convergence once some set of criteria has been met. Perhaps no improvement for X iterations.

2.2 Genetic algorithm and simulated Annealing as examples

Genetic algorithm: need to explore large portions of the parameter space at random. Concept of “neighbor” is vague.

A nice shiny app¹

An GA example: Since a new treatment for Hep C has become available, where is the optimal place to locate limited new Hep C resources, considering where our patients live?

The problem become intractable with large number of locations and resources: How many combinations of patients and clinics can I calculate the full feature space for to find a maximum?

- Exact Solution is NP-Hard
- Calculations = $n^{\sqrt{k}}$
- I conveniently stopped my analysis at 6 sites with ~5k patients, requiring 1,149,712,053 distance calculations (I have a big server)
- The “k-center” problem

Simulated Annealing:

- The concept of a ‘neighbor’ is strong.
- Can be sensitive to parameter choice, or algorithm gets stuck in global minima!
- Generally, you should try both to see what works best. Hard to guess up front.

2.3 Constrains

Hard constraints

¹<https://toddschneider.com/posts/traveling-salesman-with-simulated-annealing-r-and-shiny/>

- If this constraint is violated, we have invalid solution.
- Labor laws, number of nurses available, etc

Soft Constraints

- These are nice to meet if possible (included in cost function somehow), but if they are not met the solution is still valid.
- Nurse prefers to only work X night shifts per month.
- Leave requests.

第三章 Continuous optimization

Points to learn in continuous optimization:

- Understand first and second derivatives and the role they play in optimizing continuous functions.
- Understand general steps in continuous optimization
- Contrast 1st order versus 2nd order derivative optimization methods
- Extend these thoughts to the distributed computing context

Things to consider for smart steps:

- Initialization value: where should I start?
- Direction of the gradient: what direction should I we step towards?
- Step size: how big of a step should we take?

3.1 First and second derivative

First derivative/gradient

- Instantaneous slope of a point (rate of change of the function)
- If we have multiple input variables (multiple x 's), then we need to know the gradient in the direction of each of the x 's (partial derivatives). The matrix of partial first derivatives is called **the Jacobian**.

Second derivative

- Tells me the 'curvature' of a function.

- Rate of change of the first derivative.
- If we have multiple input variables (multiple x 's), then the matrix of partial second derivatives is called **the Hessian**.

A comparison of first and second order derivative methods

- Second order derivative methods are generally more accurate and converge in fewer steps
- Second order derivative methods are more resource intensive
- Sometimes it is easy and cheap to calculate the Hessian... (generalized linear models with canonical link), so why not?
- Sometimes it is expensive though.
- There is a tradeoff here that is context dependent.

Why not always second derivative

- It's expensive and takes more time / resources / memory.
- The Jacobian matrix only requires $O(n)$ storage.
- The Hessian matrix requires $O(n^2)$ storage.
- The size of the matrix grows exponentially with the size of the input data (specifically the number of columns).
- But... it can be more efficient if we take fewer steps, as long as the dataset isn't too big.

There is a tradeoff between the accuracy of the next step we take, and the amount of resources it takes to calculate the next step.

3.2 First Order Derivative methods

3.2.1 Gradient descent

Sourced from wikipedia¹

- Start somewhere (initial values for X)

¹https://en.wikipedia.org/wiki/Gradient_descent

- Calculate the gradient at that point
- Take a step in the correct direction based on the gradient
- Step size is a function of the gradient (larger gradient means larger step size)
- Repeat.
- Stop algorithm once it converges (within tolerance) to a single point.
- This is **expensive!** Requires a full pass over all training data at every step. . .

3.2.2 Stochastic gradient descent

(https://en.wikipedia.org/wiki/Stochastic_gradient_descent)

- Start somewhere (initial values for X)
- Randomly shuffle the data by row.
- For $i=1,2,3 \dots n$, calculate the gradient **only for the i 'th sample** (not the full dataset).
- Take a step in the correct direction based on the gradient
- Step size is a function of the gradient (larger gradient means larger step size)
- Repeat.
- Stop algorithm once it converges (within tolerance) to a single point.
- This is less costly, since you don't need a full pass over the data for every step. But it is less accurate as well. . .

3.2.3 Coordinate descent

(https://en.wikipedia.org/wiki/Coordinate_descent)

- If we have multiple X values, then we optimize them by only considering a change in a single X value at a time. The step size is based on only changing one X.
- This is useful if it is hard to calculate the gradient for all variables (the Jacobian), but easier to only work on one variable at a time.
- This is the optimal solver for regularized GLM's (elastic net regression).
 - Start somewhere (initial values for X)
 - Choose one of the X's (coordinates), change the value (your step).

- Calculate the objective function. Next round, change a different X .
- Repeat.
- Stop algorithm once it converges (within tolerance) to a single point.

Reference to Regularization Paths for Generalized Linear Models via Coordinate Descent²

3.3 Second Order Derivative methods

3.3.1 Iteratively reweighted least squares (IRLS)

- This is the most popular in data science frameworks.
- This is efficient and accurate for generalized linear models (logistic regression, Poisson regression etc...)
- The Hessian matrix (second derivative) gives us information about the uncertainty of the solution. This is where our confidence intervals and p-values come from!

3.3.2 Newton-Raphson optimization and Fisher Scoring

- More general cases of IRLS. These are identical when applied to GLM's, so you often see the terms interchanged when talking about GLM's. These are not necessarily identical outside of GLM's.
- Second-order methods are sometimes termed 'Newton methods'

3.3.3 Limited-memory Broyden-Fletcher-Goldfarb-Shanno (L-BFGS)

Technically first order since it does not evaluate the Hessian.

²<https://web.stanford.edu/~hastie/Papers/glmnet.pdf>

- However, it does approximate the Hessian by storing the prior gradient evaluations!
- So we get some idea of the rate of change of the gradient by looking at the trend of the prior gradients.
- This is termed ‘quasi-Newton’ since we approximate the Hessian without actually incurring the full cost

Orthant-Wise Limited-memory Quasi-Newton (OWL-QN) is an extension to this that effectively optimizes regularized regression (L1 or elastic net). This is implemented in Apache Spark.

3.3.4 L-BFGS Versus IRLS for GLM’s

- Both can be implemented in parallel by calculating chunks of rows at a time.
- Consider m rows and n columns. . . the IRLS algorithm requires an $N \times N$ matrix be generated no matter how small we make M by chunking by row.
- So if we have a large number of columns, IRLS can underperform (take too long / too much memory), even in distributed environments.
- L-BFGS is more efficient for a large number of columns. But, it is generally less accurate (Takes more steps).

3.4 Close thoughts

- Choice of optimization method is important!
- Depending on how large your data is, or how complex your objective function is, you may have to try different optimization methods.
- If the optimization method you choose does not give you estimates about the uncertainty of the solution (I.E. confidence intervals and p-values), you may be able to get that from a direct Hessian calculation once you have declared the optimal solution to be found.

第四章 Introduction

第五章 Introduction

第六章 Introduction

