

CS 7641 Machine Learning

Assignment 2

Philip Bale
pbale3

Due Sunday March 11th, 2018 11:59pm

Part 1: Neural Network Optimization

Introduction

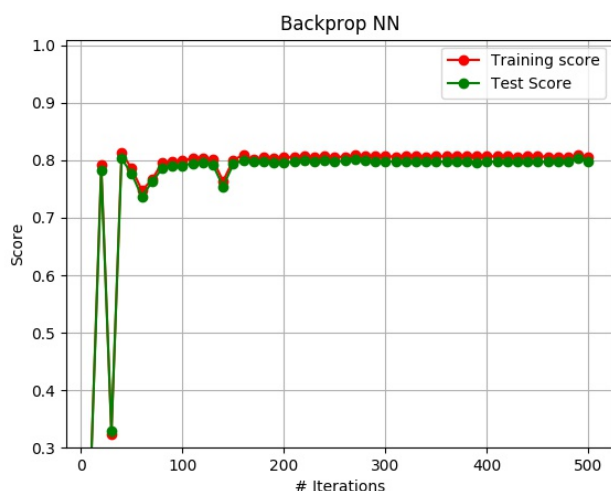
Part 1 of the assignment surrounds using randomized optimization to find the best possible weights for a specific neural network. In assignment 1, backpropagation was used to find optimal parameters for a neural network. This neural network took in various input features for US permanent visa applicants and then attempted to predict the outcome of an application. After various tests, I found the optimal parameters of: 6-node input layer, one hidden layer with 100 nodes, one output node, and about 500 iterations.

I chose this problem because, as someone who has worked with a large number of first-generation visa holders and immigrants, I am extremely interested in building tools to help others to achieve the same. At the end of the day, the goal is it to try to determine the application result before time, money, and other resources are spent.

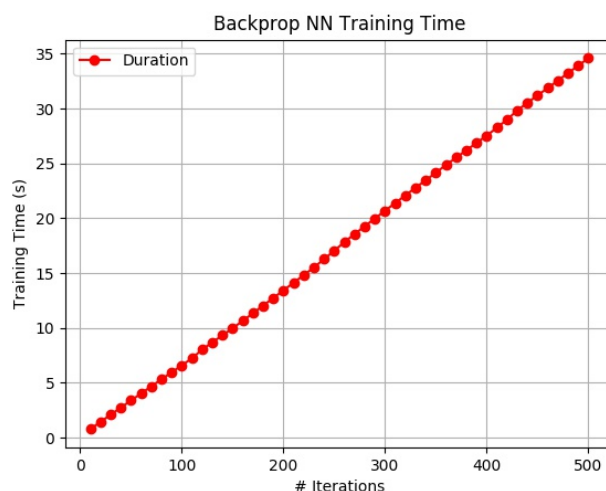
1) Backpropagation (Assignment 1)

Overview

The first weight-finding algorithm used was backpropagation. Backpropagation works by essentially calculating the error at the end of a network, and then working backwards to minimize that error over various iterations. An error (or loss) function is effectively minimized over time using this backpropagation technique. As discussed in assignment 1, the permanent visa is rather large and robust. It is quickly learnable by various different learners and in such backpropagation found significant success.



Permanent Visa Applicant NN Learning Curve



Permanent Visa Applicant NN Learning Curve

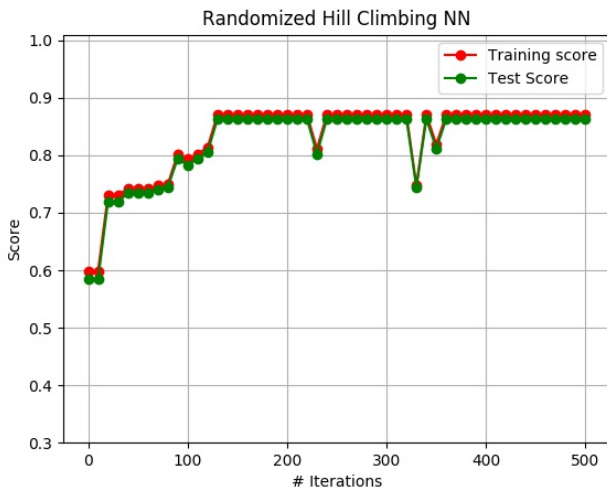
Right around 50 iterations, the network begins to converge at around an 80% success rate. Seeing as the training and test score track either rather closely, it is apparent that the dataset is rather robust and consistent. One thing to note is that the training time scales linearly with the number of iterations—which makes sense since the same amount of calculations with similar complexity are performed on each iteration of backpropagation.

2) Randomized Hill Climbiing

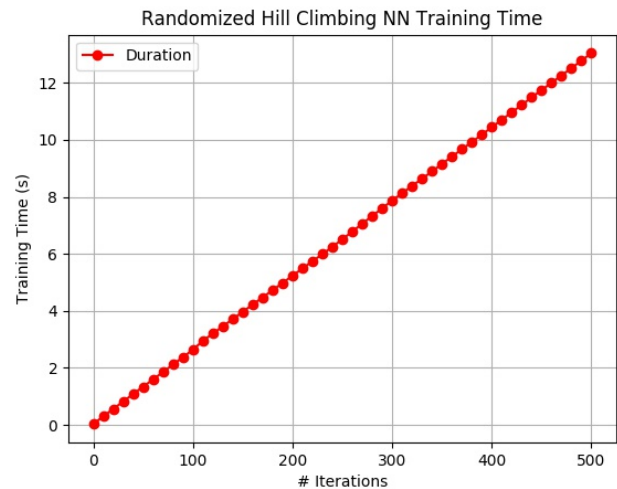
Overview

The second weight-finding algorithm used was randomized hill climbing. Randomized hill climbing works by taking a random starting point and then incrementally attempting to improve on that point. In the context of a neural network trying to find weights, randomized hill climbing selects random weights and then moves in a direction so as to try to find a better result for that weight—akin to trying to move up an optimization 'success hill'.

One thing to note is that we are using randomized hill climbing, not random restart hill climbing. In such, the algorithm is prone to getting caught in local optimizations, or local maximums.



Permanent Visa Applicant NN Learning Curve



Permanent Visa Applicant NN Learning Curve

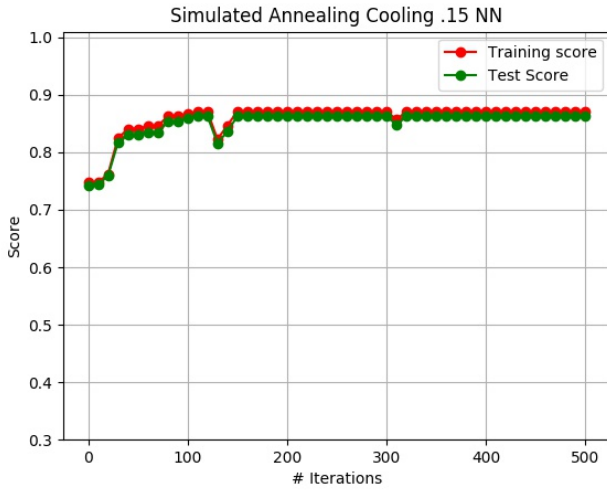
High accuracy results are achieved right around 125 iterations+. While this is subject to randomness, by looking at the results it is shown to become rather consistant. By looking at the score results, one can see various instances of the randomized hill climbing getting caught in local optimas and being unable to escape. Such is the case around 220, 320, and 350 iterations.

Similar to backpropogation, training time scales linearly with the number of iterations run. Training time tends to be a bit faster using randomized hill climbing because of a reduction in calculations necessary. Whereas backpropogation needed to do calcuations to minimize error moving backwards through the network, randomized hill climbing simply needs to move in one direction and determine if the new weights are better.

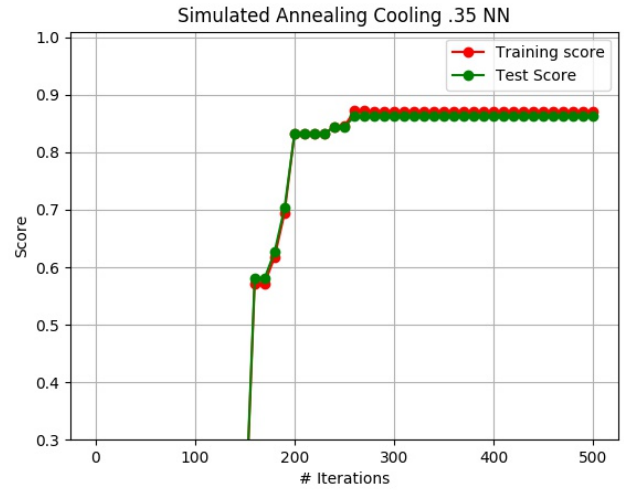
3) Simulated Annealing

Overview

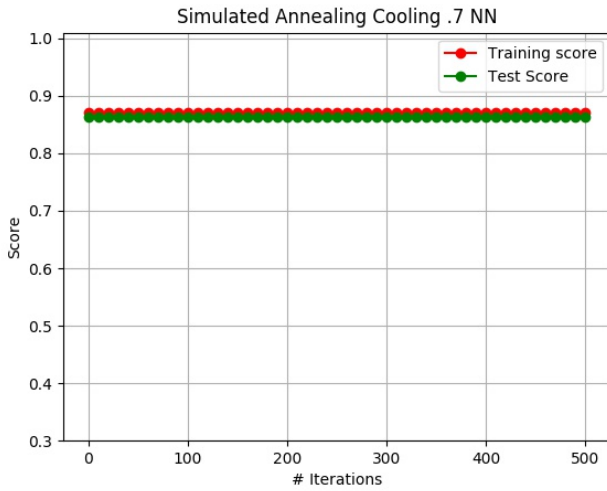
The third weight-finding algorithm used was simulated annealing. Simulated annealing works by taking a random solution and then samples nearby alternatives. By comparing the alternatives to the original solution, the optimizer decides to either stick with the original solution or move to the new one. Using a temperature and cooling parameter, the algorithm is more open to worse solutions at first but gradually moves towards only accepting better solutions.



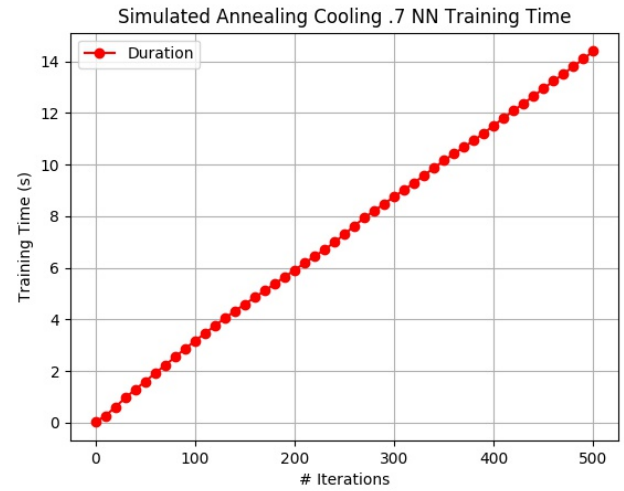
Permanent Visa Applicant NN Learning Curve



Permanent Visa Applicant NN Learning Curve



Permanent Visa Applicant NN Learning Curve



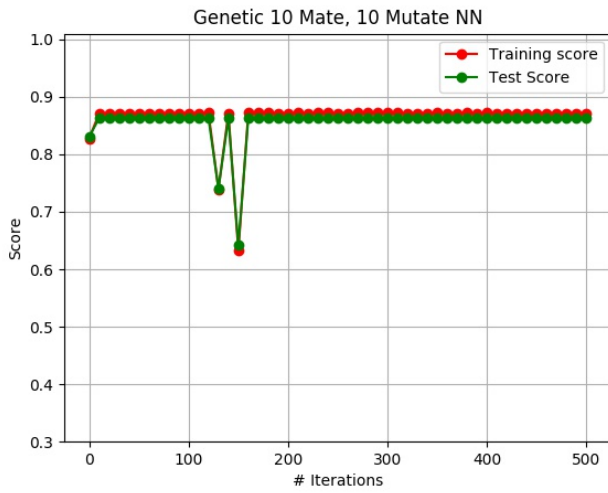
Permanent Visa Applicant NN Learning Curve

4) Genetic Algorithms

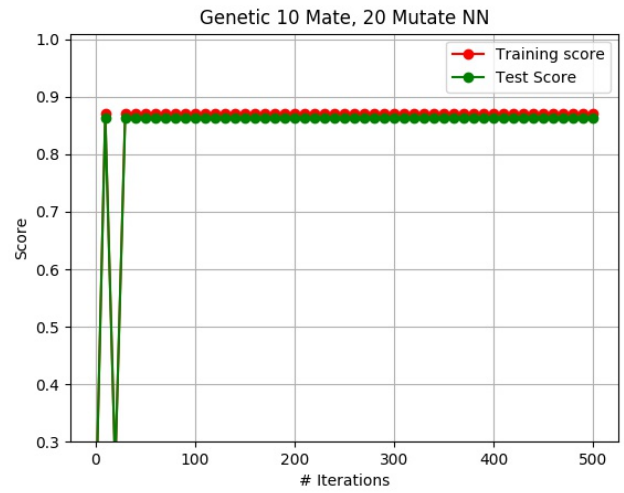
Overview

The fourth weight-finding algorithm used was a genetic algorithm. Genetic algorithms work by starting with an initial solution and then making modifications in an attempt to improve the solution. The modifications generally allowed are mutation (changing random parts of the solution), crossover (taking specific sections from various solutions and combining them), and selection (selecting certain sections from a solution to use again). In the context of a neural network, we can use genetic algorithms to make modifications to our network's weights.

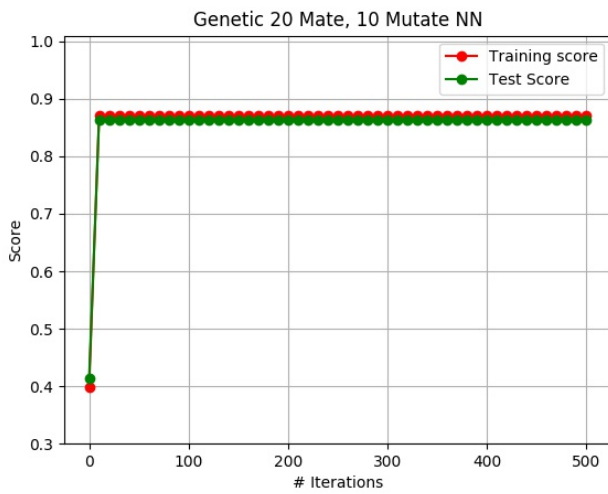
In our testing, a population size of 50 was used in order to get a diverse initial sampling of possible solutions. The number of instances to mate and to mutate was varied throughout the trials.



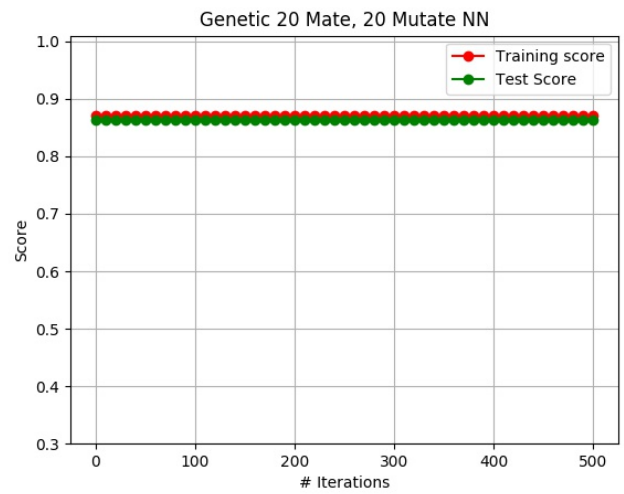
Permanent Visa Applicant NN Learning Curve



Permanent Visa Applicant NN Learning Curve

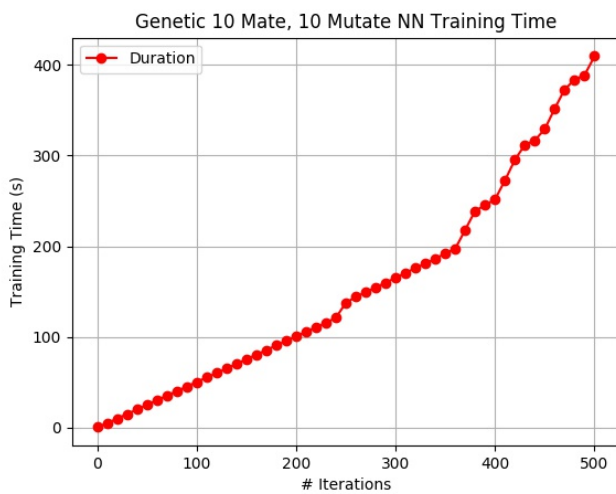


Permanent Visa Applicant NN Learning Curve

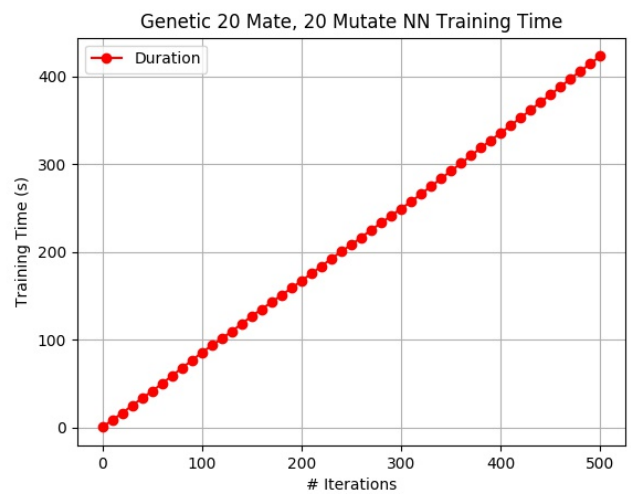


Permanent Visa Applicant NN Learning Curve

More discussion



Permanent Visa Applicant NN Learning Curve



Permanent Visa Applicant NN Learning Curve

Conclusion

The optimization algorithms used all inevitably produced similar results. What differed, however, was the amount of training time required, the tuning of parameters, and how many iterations were required to achieve a consistent, optimal solution. All-in-all, the genetic algorithm approach consistently produced the best results for our network. This makes sense, as the data was rather homogenous and there appeared to be very few outliers. By learning the training data well, the learner was able to perform similarly well on the (nearly identical) test data.

Part 2: Optimization Problem Domains

Introduction

For this part of the assignment, three different optimization problems are examined: continuous peaks, the traveling salesman, and flip flop. The three optimization techniques from above are used (randomized hill climbing, simulated annealing, and genetic algorithms) as well as another algorithm, MIMIC.

MIMIC, similar to the other algorithms, works to find the globally optimal solution. Unlike the other algorithms, however, it retains knowledge of previous iterations and uses this information to more efficiently find better solutions. MIMIC is particularly strong in regards to problems that maintain patterns between subsets of their parameters.

1) Continuous Peaks

Overview

The continuous peaks problem is an extension of the four peaks problem that allows for a wide variety of local maximums. This algorithm is particularly interesting due to the potentially large number of local maximums. It is especially difficult for an algorithm that cannot escape local peaks to perform well.

2) Traveling Salesman

Overview

The traveling salesman problem is a commonly used, old problem that focuses on a salesperson trying to minimize their round-trip distance between any number of cities. This problem is particularly interesting due to its real-world implications, such as route planning for delivery trucks and everyday errands—it also has no known polynomial time solution!

3) Flip flop

Overview

The flipflop problem is another common optimization, where one attempts to count the number of bits that alternate with its next neighbor in a bit string. This problem is particularly interesting because, since the strings are randomized, there is significant potential for a large number of local minimum and maximums.

Conclusion

asdf

- Job features: Industry code, job class, wage rate, wage type
- Geographic features: Country of citizen and employer location

The classes observed for this dataset are simply 'approved' and 'denied'. The dataset contains 374365 total samples.

Why is the dataset interesting?

This dataset is interesting due to its potential to aid in the visa application process from a cost and time savings potential. It could also enable confidence in those interested in applying for a US permanent visa but doubting their chances of acceptance. At the end of the day, the goal is it to try to determine the application result before time, money, and other resources are spent. As someone who has worked with a large number of first-generation visa holders and immigrants, I am extremely interested in building tools to help others to achieve the same.

From a machine learning perspective, the dataset is incredibly interesting due to its wide variety of features and the variety of values those features can take. An immense number of job types, wage rates, and citizenships alone create an

extremely diverse dataset. Additionally, the number of samples available provide a comprehensive picture of historical data, lending towards greater confidence in training and testing rates.

2) Home Sale Price Predictions

Overview

The second classification problem revolves around classifying a home's price bracket based upon the various characteristics of the home. Among the features used in the classification are :

- Subjective measurements: Exterior condition, house style, overall quality rating, and overall condition
- Objective measurements: Type of dwelling, building type, lot size, neighborhood, year built, and year sold

After an initial review of the dataset, the classes were defined as pricing brackets divided into 100k groups. I.e: 0-100k, 100k-200k, 200k-300k, etc. The dataset contains 1451 samples. An additional dataset containing another 1400 testing samples exists but was not used as it contains unclassified sale prices. It will, however, prove useful for unsupervised learning.

Why is the dataset interesting?

This dataset is interesting for two primary reasons: real-world applicability and participating in a Kaggle challenge. First, modeling home prices is both a difficult and lucrative task. If one can successfully model home sale prices on large sets of data, he/she can make large amounts of money investing in real estate when he/she detects outliers in listed price vs. what it is expected to sell for. This applies to flipping, investing, and remodeling. Second, the dataset is part of an ongoing Kaggle competition that does not have a winning solution yet. By taking part of the competition, the dataset presents the opportunity to work towards a winning solution and advance ones algorithms over time.

Houses can have a very large amount of features—with a large amount of variety in the individual features. Similarly, housing is prone to personal taste and frequent need for upgrades/modernization. In such, I believe price estimation is an excellent problem, full of depth and complexity, that is suitable for a machine learning approach.

General Data Processing

The datasets I used were both relatively clean to begin with. One small problem, however, was that a lot of my features on both datasets were text-based. To transform the features into numeric values suitable for the machine learning algorithms, I used a label encoder built into ScikitLearn.

I also did a small amount of preprocessing of the data to make it more suitable for classification. I dropped all unnecessary columns to help speed up with data processing in general—which proved immensely helpful when dealing with the more computationally intensive algorithms. In the case of home prices, I precalculated the brackets based on the 'sale price' data label. In the case of visa applications, I pregrouped the case outcome so that results such as 'certified' and 'certified withdrawn' are both concerned as 'approved' conditions whereas 'denied', 'invalidated', and 'rejected' all resolve to 'denied'.

1: Decision Trees

A decision tree classifier was the first algorithm applied to the datasets. Various values of max_depth were tested as a means of pruning unnecessary leaves. Similarly, a grid search was used to test whether a 'gini' or 'entropy' criterion was more effective.