

MAT342 NUMERICAL ANALYSIS FINAL PROJECT

**OPTIMIZATION OF MACHINE LEARNING
CLASSIFIERS WITH DIRECT SEARCH
ALGORITHMS**

Rachel You
Gordon College
April 28, 2017

1 Introduction

In machine learning, besides good choices of classifiers and abundant data, good parameters for the classifier chosen are also essential in determining the goodness of classification. Finding the best parameters is an optimization problem. Optimization problems can also all be transformed into minimization problems, and be resolved with minimization algorithms. Minimizing functions in a multivariate case can be complicated when we are not able to take the derivative of the function. Nelder-Mead Algorithm and Method of Simulated Annealing are two methods to minimize a function without taking the derivative. We are going to examine each algorithm and apply each to tuning Support Vector Machine (SVM).

2 Nelder-Mead Algorithm

Nelder-Mead Algorithm starts with an initial simplex, and proceeds by iterations of operations of reflection, expansion, contraction, and shrinkage. Eventually, the simplex will reach the optimal solution.

2.1 Algorithm

Suppose we are solving a problem with n -simplex. At each iteration, we use $n+1$ points, denoted by $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_{n+1}$, with $f(\mathbf{x}_1) \leq f(\mathbf{x}_2) \leq \dots \leq f(\mathbf{x}_{n+1})$. A simplex S_k is denoted by $S_k = \langle x_1, x_2, \dots, x_{n+1} \rangle$. In two-dimensional, it will be 3 points, forming a 3-simplex, a triangle.

Trial steps are generated by the operations of reflection, expansion, contraction, and shrinkage. A reflected vertex is computed by reflecting the worst vertex, \mathbf{x}_{n+1} , through the centroid of the remaining vertices.

$$x_r = (1 + \alpha)\bar{x} - \alpha\mathbf{x}_{n+1},$$

where $\alpha = 1$, and \bar{x} is the centroid defined by

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i.$$

The reflected vertex is accepted if $f(\mathbf{x}_1) \leq f(x_r) \leq f(\mathbf{x}_n)$.

If the $f(\mathbf{x}_r) \leq f(\mathbf{x}_1)$, then we produce an expansion.

$$x_e = \gamma x_r + (1 - \gamma)\bar{x},$$

where $\gamma = 2$. The expansion vertex is accepted if $f(x_e) < f(\mathbf{x}_1)$, otherwise the reflected vertex is accepted.

If $f(\mathbf{x}_n) \leq f(x_r)$, then a contraction is computed. If $f(\mathbf{x}_{n+1}) \leq f(x_r)$, then the internal contraction vertex is computed as

$$x_c = \beta \mathbf{x}_{n+1} + (1 - \beta)\bar{x},$$

otherwise, the external contraction vertex is computed as

$$\hat{x}_c = \beta x_r + (1 - \beta)\bar{x},$$

where $\beta = \frac{1}{2}$. The contraction vertex is accepted if it has a lower function value than \mathbf{x}_n .

If both reflection vertex and contraction vertex are rejected, then the simplex is shrunk. Each vertex \mathbf{x}_i , except the best point \mathbf{x}_1 , is replaced by the point halfway between \mathbf{x}_i and \mathbf{x}_1 .

$$x_i \leftarrow \frac{\mathbf{x}_i + \mathbf{x}_1}{2}.$$

Finally, function values of the accepted points are sorted with the remaining point(s), and the next iteration begins.

2.2 Stopping Criteria

There are at least multiple stopping criteria. Here we introduce four, the first three discussed by Dennis and Woods, and the last one discussed by Cheney and Kincaid.

The first one is to halt when the standard error of the function values falls below some threshold value:

$$\frac{1}{n} \sum_{i=1}^{n+1} (f(x_i) - \bar{f})^2 < \epsilon_1,$$

where \bar{f} is the average of the function values and $\epsilon_1 > 0$ is a preset tolerance value.

Another stopping criterion bases on how far the simplex moves at an iteration, halting when:

$$\frac{1}{n} \sum_{i=1}^n \|\mathbf{x}_i^k - x_i^{k+1}\|^2 < \epsilon_2,$$

where $\epsilon_2 > 0$ and x_i^{k+1} is the i^{th} unordered point in the $k + 1^{\text{st}}$ simplex. The main objection to this method is that the left-hand side for a shrinkage step will be greater than the value for a contraction step, and shrinkage occurs frequently when the simplex is in a neighborhood of a local minimizer.

Therefore, the third stopping criterion is introduced by Woods:

$$\frac{1}{\Delta} \max_{2 \leq i \leq n+1} \|\mathbf{x}_i - \mathbf{x}_1\| \leq \epsilon_3,$$

where $\Delta = \max(1, ||(x)_1)$ and $\epsilon > 0$. This is a measure of the relative size of the simplex.[1]

The fourth stopping criterion tests whether the relative flatness is small[2], which is

$$\frac{F(x_0) - F(x_n)}{|F(x_0)| - |F(x_n)|} < \epsilon_4$$

2.3 Advantages

- Robustness: It tolerates noise in the function values.
- Simplicity in programming: Trail points are obtained using very simple algebraic manipulations and these points are accepted or rejected based only on their function values.
- Low overhead in storage and computation: When the number of variable is small, this algorithm is often competitive with much more complex algorithms that require a great deal of overhead in storage and algebraic manipulations.

2.4 Disadvantage

- Stopping criteria is not guaranteed to converge.

3 Method of Simulated Annealing

Method of Simulated Annealing (SA) is good at minimizing difficult functions, especially discrete functions. It uses probabilities to assign new values to each data point and evaluate to get to another iteration.

3.1 Algorithm

The algorithm generates a sequence of points x_1, x_2, x_3, \dots and it's hoped the the minimum of function values of the points will converge.

For iteration k, we first generate random points u_1, u_2, \dots, u_m in a large neighborhood of x_k . The minimum function value will be selected:

$$F(u_j) = \min\{F(u_1), F(u_2), \dots, F(u_m)\}$$

If the newly computed function value is less than function value of x_k , then set $x_{k+1} = u_j$. If $F(u_j)$ is not better than $F(x_k)$, we assign a probability p_i to u_i by formula:

$$p_i = e^{\alpha[F(x_k) - F(u_i)]} (1 \leq i \leq m)$$

where α is a positive parameter set by user code. The probabilities are normalized by dividing each by their sum.

$$S = \sum_{i=1}^m p_i$$

$$p_i \leftarrow p_i / S$$

Finally, a point is chosen randomly from points u_1, u_1, \dots, u_m with probabilities p_i taking into account.[2]

3.2 Advantages

- Avoid getting stuck at local optimums: the complicated choice of x_{k+1} can help considering points that are not stuck at the same local minimum.[3]

3.3 Disadvantages

- Repeatedly annealing with a $1/\log(k)$ schedule is very slow, especially if the cost function is expensive to compute.
- For problems where the energy landscape is smooth, or there are few local minima, SA is overkill — simpler, faster methods (e.g., gradient descent) will work better. But generally don't know what the energy landscape is for a particular problem.
- Heuristic methods, which are problem-specific or take advantage of extra information about the system, will often be better than general methods, although SA is often comparable to heuristics.
- The method cannot tell whether it has found an optimal solution. Some other complimentary method (e.g. branch and bound) is required to do this.[4]

4 Python Implementations and Simple Examples

Here a small example is implemented to visualize both methods in 2D.

The function we are using is $f(x, y) = x^2 + y^2$, for which we can easily know the minimum to be at $(0, 0)$.

One thing to be noted is that although the function used here has 2 variables, the implementation is generalized and can be used for higher dimensions.

4.1 Nelder-Mead Implementation with 2D Example

Here we use the third stopping criterion, which does not require computation of the function values. It would be efficient when applied in selecting machine learning parameters since having the classifier fit the training data usually takes a long time.

```
import numpy as np
import matplotlib.pyplot as plt

# define the function
def f(x):
    return np.power(x[0],2) + np.power(x[1],2)

# define constants to be used
alpha = 1.0
gamma = 2.0
beta = 0.5
epsilon = 0.01

# initialize x array and other variables
x = np.array([[100.0,150.0],[120.0,-90.0],[-110.0,-130.0]])
fx = np.array([f(x_) for x_ in x])
fxsort = fx.argsort()
fx = fx[fxsort]
x = x[fxsort]
n = 2
count = 0

# draw contour plot
xlist = np.linspace(-1.0, 1.0, 100) # Create 1-D arrays for x,y dimensions
ylist = np.linspace(-1.0, 1.0, 100)
X,Y = np.meshgrid(xlist, ylist) # Create 2-D grid xlist,ylist values
Z = f(np.array([X,Y]))
plt.contour(X, Y, Z, [10000.0, 20000.0, 30000.0, 40000.0], colors = 'b',
            linestyle = 'solid')

# iteration
# reflection
while True:
    count += 1
    xnew = []
    xbar = 1/float(n)*np.sum(x[:n-1,:],axis=0) # centroid
```

```

xr = (1+alpha)*xbar - alpha*x[n]
fxr = f(xr)
if fx[0] <= fxr <= fx[n-1]:
    xnew = xr #reflection_accepted = true
    print "reflection"
elif fxr <= fx[0]:
    # expansion
    xe = gamma*xr + (1-gamma)*xbar
    if f(xe) < fx[0]:
        xnew = xe # expansion_accepted = true
        print "expansion"
    else:
        xnew = xr # reflection_accepted = true
        print "reflection"
else: # fx[n-1] <= fxr:
    # contraction
    if fx[n] <= fxr:
        # internal contraction
        xc = beta*fx[n]+(1-beta)*xbar
    else:
        # external contraction
        xc = beta*xr + (1-beta)*xbar
    if f(xc) < fx[n-1]:
        xnew = xc
        print "contraction"
    else: # both reflection vertex and contraction vertex are rejected
    # shrinkage
    for i in range(1,n):
        x[i] = (x[i] + x[0])/2.0
    xnew = (x[n] + x[0])/2.0
    print "shrinkage"
x[n] = xnew
# resort the array
fx = np.array([f(x_) for x_ in x])
fxsort = fx.argsort()
fx = fx[fxsort]
x = x[fxsort]
# plot the simplex
p = plt.Polygon(x, closed=True, fill=False)
ax = plt.gca()
ax.add_patch(p)

```

```

Delta = max(1.0,np.sum(np.abs(x[0])**2,axis=-1)**(1./2))
norm_array = np.array([np.sum(np.abs(x[i] - x[0])**2,axis=-1)**(1./2)
                        for i in range(1,n+1)])
rel_size = 1.0/Delta*max(norm_array)
if rel_size < epsilon:
    break

bestX = x[0]
bestF = fx[0]

print "number_of_Iterations:"
print count
print "Best_variable_result:"
print bestX
print "Best_function_value:"
print bestF

annotation = "f(%3.8f,%3.8f)=%3.8f" % (bestX[0],bestX[1],bestF)
# plot
plt.axis([-70,150,-150,40])
plt.title("Nelder-Mead_Simple_2D_Example_with_f_=_x^2_+_y^2")
plt.xlabel("x")
plt.ylabel("y")
plt.plot(bestX[0],bestX[1],'^')
plt.annotate(annotation, (bestX[0],bestX[1]))
plt.show()

```

Output shows the operations performed in each iteration and the total number of iterations taken. (Plot: Figure 1)

```

shrinkage
contraction
shrinkage
shrinkage
reflection
shrinkage
shrinkage
contraction
shrinkage
shrinkage
contraction
shrinkage

```

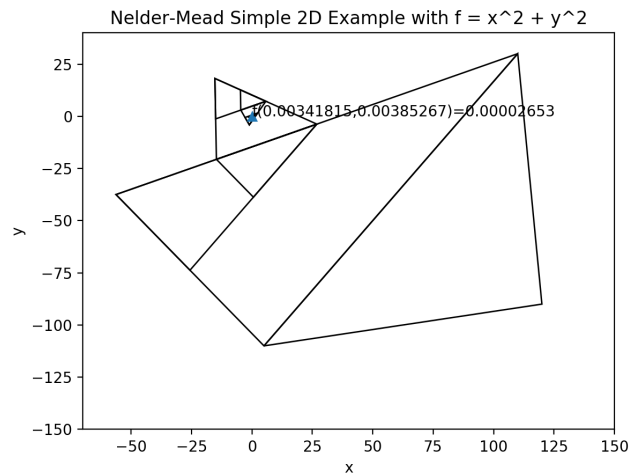



Figure 1: Nelder Mead Example

```

contraction
contraction
contraction
reflection
reflection
reflection
contraction
contraction
contraction
contraction
number of iterations:
22
Best variable result:
[ 0.00341815  0.00385267]
Best function value:
2.65268148309e-05

```

4.2 Simulated Annealing Implementation with 2D Example

Since Simulated Annealing does not have a good stopping criterion, multiple trials have to be performed to find a good number of iterations.

```

import numpy as np
import matplotlib.pyplot as plt

```

```

# define the function
def f(x):
    return np.power(x[0],2) + np.power(x[1],2)

# define constants to be used
n = 2
iterations = 22
m = 100
radius = 15.0
sigma = 5.0
alpha = 1.0

# define initial conditions
x = np.zeros((iterations,2))
fx = np.zeros(iterations)
x[0] = [100.0,150.0]
fx[0] = f(x[0])

#iteration
for k in range(0,iterations-1):
    u = np.zeros((m,n))
    count = 0
    while count < m:
        unew = np.zeros(n)
        for i in range(0,n):
            unew[i] = np.random.normal(x[k,i],sigma)
        distance = np.sum(np.abs(unew-x[k])**2,axis=-1)**(1./2)
        if distance < radius:
            u[count] = unew
            count += 1

    fu = np.array([f(u_) for u_ in u])
    j = np.argmin(fu)

    # accept the new variable if function value gets better
    if fu[j] < fx[k]:
        x[k+1] = u[j]
        fx[k+1] = f(x[k+1])
    else:
        p = np.zeros(m)

```

```

    for i in range(0,m):
        p[i] = np.exp(alpha*(fx[k]-fu[i]))
    S = np.sum(p)
    p = p/S
    xi = np.random.rand()
    for i in range(0,m):
        if xi < np.sum(p[:i]):
            x[k+1] = u[i]
            fx[k+1] = f(x[k+1])

bestX = x[np.argmin(fx)]
bestF = np.min(fx)

# print the function values changing
print fx
print "Best_variable_result:"
print bestX
print "Best_function_value:"
print bestF

annotation = "f(%3.8f,%3.8f)=%3.8f" % (bestX[0],bestX[1],bestF)
# plot
plt.plot(x[:,0],x[:,1])
plt.title("Simulated_Annealing_Simple_2D_Example_with_f=_x^2+_y^2")
plt.xlabel("x")
plt.ylabel("y")
plt.plot(bestX[0],bestX[1],'^')
plt.annotate(annotation, (bestX[0],bestX[1]))
plt.show()

```

Output: (Plot: Figure 2)

```

[ 3.25000000e+04  2.90169635e+04  2.49033265e+04  2.10309765e+04
 1.81313495e+04  1.48308788e+04  1.22384950e+04  9.89537276e+03
 7.61079384e+03  5.48629955e+03  3.73471473e+03  2.37701988e+03
 1.46480064e+03  7.09729534e+02  1.83382094e+02  1.91799304e-01
 6.19211705e+01  2.67922432e-01  1.64130436e+02  1.79173840e+01
 1.15513015e-01  8.04819945e-02]
Best variable result:
[-0.14821482  0.24189742]
Best function value:
0.0804819945256

```

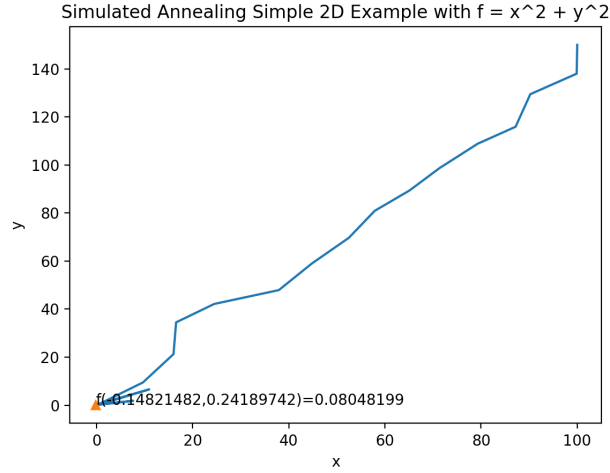


Figure 2: Simulated Annealing Example

4.3 Comparison

With same amount of steps, Nelder-Mead Algorithm gets to really small error. However, what Simulated Annealing gets is not that accurate. In fact, we are not able to control the accuracy of Simulated Annealing since it's using random number generator.

Also, Nelder-Mead Algorithm converges much faster, while Simulated Annealing doesn't converge.

In other functions, we might be able to test out the advantage of Simulated Annealing of avoiding getting stuck at a local minimum.

Method	steps	[x y]	f
Nelder-Mead	22	[0.00341815 0.00385267]	2.65268148309e-05
Simulated Annealing	22	[-0.14821482 0.24189742]	0.0804819945256

Table 1: An example table.

5 Application to SVM

Support Vector Machine (SVM) is a machine learning algorithm. Two essential numerical parameters are C and gamma. C controls tradeoff between smooth decision boundary and classifying training points correctly, and larger C tends to make training points more correct. Gamma controls whether or not the boundary is smooth. Here, we are going to use both Nelder-Mead Algorithm and Simulated Annealing to find the best parameters for

an SVM with python package "sklearn" that recognizes hand written numerical digits that are represented in 28×28 -pixel arrays.

Since the methods originally are designed to find minimums, we take the negatives of accuracy scores to turn the maximization problem into a minimization problem.

5.1 Nelder-Mead

With Nelder-Mead method, it seems that starting with random initial conditions will lead local optimum.

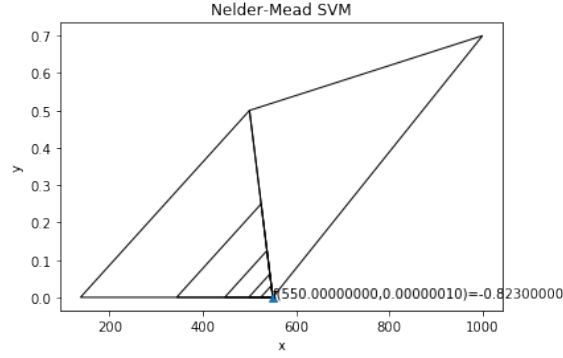


Figure 3: I.C. $x = [[100, 1], [500, 0.5], [1000, 0.7]]$, epsilon = 0.05

When we start with reasonable initial condition $x = [[1000, 1], [5000, 0.1], [10000, 0.01]]$, it in the end finds the best fit at $[10000, 0.01]$ with accuracy score 0.96433333 after 8 rounds.

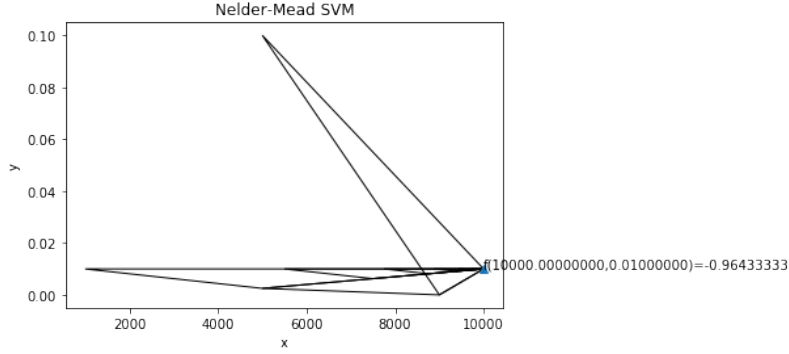


Figure 4: I.C. $x = [[1000, 1], [5000, 0.1], [10000, 0.01]]$, epsilon = 0.05

In fact, Nelder-Mead depends so much on initial condition in this case especially when we don't use a very large epsilon value. If we change the initial condition to $x = [[10002, 1.1], [5050, 0.09], [9999, 0.02]]$, it gives a better accuracy score 0.967 at $[9999, 0.02]$.

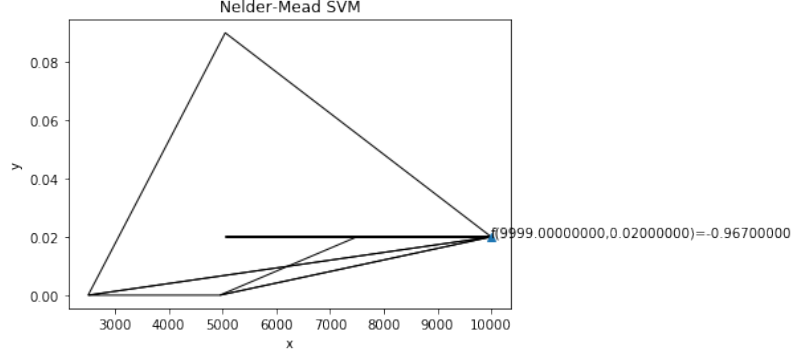


Figure 5: I.C. $x = [[10002, 1.1], [5050, 0.09], [9999, 0.02]]$, $\epsilon = 0.05$

5.2 Simulated Annealing

Simulated Annealing is taking really long time when having the classifier to fit the data, since it generates a lot of random points each iteration. Randomly selected initial value is also fatal and makes it not knowing which direction to go.

When the initial parameters is set to a ridiculous point $x = [100.0, 150.0]$, the accuracy score never goes to more than 0.10833333.

If we use a reasonable value $x = [10000.0, 0.1]$ as initial condition, generating 5 iterations with each iteration taking 10 random points, we are able to get an accuracy score of 0.9676666667.

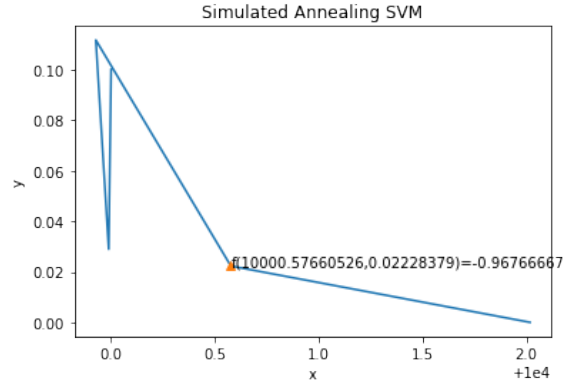


Figure 6: I.C. $x = [10000.0, 0.1]$, 5 iterations with 10 random points each iteration

5.3 Comparison to GridSearch

Comparing to sklearn's inherit GridSearch algorithm that selects the best fit from a set of values, Nelder-Mead might be able to find a better fit given a smaller epsilon, since the values to be tested are not constrained by the few choices provided. However, the initial conditions still have to be reasonable.

6 Conclusion

Nelder-Mead Algorithm and Simulated Annealing are useful for finding optimized values for functions that we are not able to calculate the derivatives. Each still has its own disadvantages. Nelder-Mead Algorithm can suffer from getting stuck at local minimums, while Simulated Annealing sometimes require too much computation resources when calculating complicated functions.

Better implementations might be available to reduce repetitive calculations of functions and thus reduce time needed. Rescaling functions might also allow the algorithms to find better results when each axes are not of the same scale.

In selecting parameters for SVM, it might be better to first make a simplex that expands across a large area, finding the relatively best point with Nelder-Mead Algorithm, and then use Simulated Annealing at that specific area to find better results. This may diminish the advantage of Simulated Annealing for not getting stuck at local optimums, but we have to make the sacrifice in regard to the overly high computational expensiveness of fitting.

References

- [1] Dennis, John E. JR. and Woods, Daniel J. Optimization on Microcomputers: The Nelder-Mead Simplex Algorithm. *New Computing Environments: Microcomputers in Large-scale Computing*, 116-121.
- [2] Cheney, Ward. and Kincaid, David *Numerical Mathematics and Computing*, 580-582.
- [3] Simulated Annealing
<http://bamboo.ee.ntu.edu.tw/LabWebsite/Frame/D0/Simulated%20annealing.pdf>
- [4] Disadvantages of Simulated Annealing
<https://cs.adelaide.edu.au/~paulc/teaching/montecarlo/node140.html>