

Ridge regression in Python

Cristina Aguilera, Jesús M. Antoñanzas

Mathematical Optimization, Data Science and Engineering
UPC

May 2019

Contents

1 Objectives 2

2 Preliminary algebra 3

2.1 Objective function 3

2.2 Constraint function 3

3 Results 5

3.1 Basic analysis 5

3.2 KKT conditions 6

3.3 Full results 7

4 Code 8

1 Objectives

We have studied and used ridge - or Tikhonov regularized - regression in various instances during the course and using different programming environments (i.e. MATLAB, AMPL and R). In those cases, it was similar to a black box: just a couple of lines of code and all was done. This time, we have written a more transparent Ridge regression implementation, which uses functions in order to compute the derivatives and Hessians of both the objective function and its constraints. In spite of this, the optimization algorithm used is still unknown to us as we have not studied it quite yet. Then, we have studied the optimality of the found point and compared it to what other implementations give as a solution (namely, the AMPL implementation).

2 Preliminary algebra

Before implementing a solution, we first need to know what the analytical expression of both the derivatives and the Hessians associated to our objective function and its constraints (in this case, just one) are.

2.1 Objective function

The objective function is

$$f(w, \gamma) = \frac{1}{2}(Aw + \gamma e - y)^T(Aw + \gamma e - y)$$

where $A \in \mathcal{R}^{m \times n}$, $w \in \mathcal{R}^{n \times 1}$, $\gamma \in \mathcal{R}$, $e \in \mathcal{R}^{m \times 1}$, $y \in \mathcal{R}^{m \times 1}$, e being a vector of ones. The gradient of $f(w, \gamma) \in \mathcal{R}^{(n+1) \times 1}$ is given by

$$\nabla f(w, \gamma) = \begin{bmatrix} \frac{\partial f(w, \gamma)}{\partial w} \\ \frac{\partial f(w, \gamma)}{\partial \gamma} \end{bmatrix}$$

where

$$\frac{\partial f(w, \gamma)}{\partial w} = A^T(Aw + \gamma e - y)$$

and

$$\frac{\partial f(w, \gamma)}{\partial \gamma} = w^T A^T e + n\gamma - e^T y$$

Its Hessian $\nabla^2 f(w, \gamma) \in \mathcal{R}^{(n+1) \times (n+1)}$ is given by the expression

$$\nabla^2 f(w, \gamma) = \begin{bmatrix} \frac{\partial f(w, \gamma)}{\partial^2 w} & \frac{\partial f(w, \gamma)}{\partial w \partial \gamma} \\ \frac{\partial f(w, \gamma)}{\partial \gamma \partial w} & \frac{\partial f(w, \gamma)}{\partial^2 \gamma} \end{bmatrix}$$

where

$$\frac{\partial f(w, \gamma)}{\partial^2 w} = A^T A, \quad \frac{\partial f(w, \gamma)}{\partial w \partial \gamma} = A^T e, \quad \frac{\partial f(w, \gamma)}{\partial \gamma \partial w} = (A^T e)^T, \quad \frac{\partial f(w, \gamma)}{\partial^2 \gamma} = n \quad (1)$$

2.2 Constraint function

The constraint function is

$$h(w, \gamma) = \|w\|_2^2 \leq t$$

Its gradient $\nabla h(w, \gamma) \in \mathcal{R}^{(n+1) \times 1}$ is

$$\nabla h(w, \gamma) = \begin{bmatrix} \frac{\partial h(w, \gamma)}{\partial w} \\ \frac{\partial h(w, \gamma)}{\partial \gamma} \end{bmatrix}$$

where

$$\frac{\partial h(w, \gamma)}{\partial w} = 2w, \quad \frac{\partial h(w, \gamma)}{\partial \gamma} = 0 \quad (2)$$

and its Hessian $\nabla^2 h(w, \gamma) \in \mathcal{R}^{(n+1) \times (n+1)}$ is

$$\nabla^2 h(w, \gamma) = \begin{bmatrix} \frac{\partial h(w, \gamma)}{\partial^2 w} & \frac{\partial h(w, \gamma)}{\partial w \partial \gamma} \\ \frac{\partial h(w, \gamma)}{\partial \gamma \partial w} & \frac{\partial h(w, \gamma)}{\partial^2 \gamma} \end{bmatrix}$$

Note that in the Hessian of the constraint we have taken γ into account, and that is because in the implementation we need all the parameters (w and γ) in a single vector. This procedure results in a matrix with the only elements different from 0 are 2's in the diagonal, except for the element corresponding to the position $[i, i]$, i being the position of γ in the aforementioned parameter vector. For us, γ was the last in the vector ($i = n+1$).

$$\nabla^2 h(w, \gamma) = \begin{bmatrix} 2 & 0 & 0 & \dots & 0 \\ 0 & 2 & 0 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & 0 & \dots & 2 & 0 \\ 0 & 0 & \dots & 0 & 0 \end{bmatrix}$$

3 Results

3.1 Basic analysis

For this data set, we have a total of 15 explanatory variables and 60 observations. The response variable is the death rate of a certain country. We first load this data set and its corresponding ridge regression model in AMPL and solve it, returning the following solution.

```
1 MINOS 5.51: optimal solution found.
2 133 iterations , objective 61484.65464
3 Nonlin evals: obj = 361, grad = 360, constrs = 361, Jac = 360.
4 w [*] :=
5 1 0.48824
6 2 -0.0463881
7 3 0.102978
8 4 -0.0377061
9 5 0.00490736
10 6 -0.0339172
11 7 -0.255786
12 8 0.00564845
13 9 0.649708
14 10 -0.12622
15 11 0.213407
16 12 -0.207831
17 13 0.109842
18 14 0.376641
19 15 0.00995978
20 ;
21
22 gamma = 895.141
23 norm2_w = -10785.1
```

Listing 1: AMPL solution of the "Death rate" regression.

As we can see, 133 iterations were needed for the routines to find a minimizer. Executing our code, the exact same solution is found in the same number of iterations.

```
1 In [0]: solution.nfev
2 Out[0]: 133
3
4 In [1]: solution.x
5 Out[1]:
6 array([ 4.88239571e-01, -4.63880870e-02,  1.02977785e-01, -3.77061350e-02,
7         4.90736137e-03, -3.39172227e-02, -2.55785705e-01,  5.64845151e-03,
8         6.49707817e-01, -1.26220227e-01,  2.13407411e-01, -2.07831309e-01,
9         1.09842480e-01,  3.76641023e-01,  9.95977558e-03,  8.95141390e+02])
```

Listing 2: Python solution of the "Death rate" regression.

Where, in our case, γ is the last value in the x vector and ω is the rest of the vector. As we can see, the value of the objective function at our solution is the same as in the AMPL one.

```
1 In [2]: solution.fun
2 Out[2]: 61484.65464046297
```

Listing 3: Value of the objective function at the found minimizer.

On the other hand, execution time with Python is not faster, as we can see taking a look at the "_solve_time" environmental variable, on which AMPL saves the "system CPU seconds used by most recent solve command" plus the "user CPU seconds used by most recent solve command".

```
1 ampl: display _solve_time;
2 _solve_time = 0.007453
```

Listing 4: AMPL solve time for the "Death rate" data set.

While in Python:

```
1 In [3]: solution.execution_time
2 Out[4]: 0.2936720848083496
```

Listing 5: Python execution time for the "Death rate" data set.

AMPL is much faster. At the end of the section we explain some of the reasons why this could happen.

3.2 KKT conditions

Knowing that our solution is the same as the AMPL one is a good indication that our point is optimal, but we can also deduce that from the output of the Python routines. Specifically, as Ridge regression is a strictly convex problem, if our minimizer meets the KKT conditions, then that it is a sufficient condition so that it is optimal, as it will always satisfy the sufficient second-order conditions. The first KKT condition is **feasibility**. Let's check that $h(\omega, \gamma) \leq t$.

```
1 In [4]: cons(solution.x)
2 Out[4]: 0.9999999999979048
```

Listing 6: Feasibility KKT condition.

The first condition is met. Now, the **second condition**: $\nabla L(\omega^*, \gamma^*, \lambda^*) = 0$.

```
1 In [5]: solution.lagrangian_grad
2 Out[5]:
3 array([-6.64695108e-08,  5.67613370e-08,  4.70404302e-08,  6.35918695e-09,
4         5.68779939e-08, -8.96648089e-10,  3.83806764e-10,  6.66005917e-09,
5         4.45743353e-08, -6.24581844e-08,  4.29736247e-09, -4.31409717e-08,
6         -1.83293196e-07,  1.04246283e-08, -2.83398265e-08, -9.45874490e-11])
```

Listing 7: 2nd KKT condition.

We can say this one is also met. Lastly, **complementary slackness**. We need for $\mu^* \geq 0$ (in this case, $\mu^* = \nu^*$) and $\nu^* g(\omega^*, \gamma^*) = 0$, where $g(\omega, \gamma) = \|\omega\|_2^2 - t$. Running a quick check:

```
1 In [6]: solution.v[0]
2 Out[6]: array([10785.1067161])
3
4 In [7]: solution.v[0]*(cons(solution.x)-t)
5 Out[7]: array([-2.25970946e-08])
```

Listing 8: Complementary slackness KKT condition.

So, KKT conditions are met at the found minimizer, meaning that, in fact, the found point is a global minimum.

A comment on speed. Our implementation is not faster than the AMPL one in spite of ours not needing to compute the gradients and Hessians symbolically. Of course, AMPL is a professional software, many times tested and improved, but we can take a guess at the reasons for this being the case are. First, because of the way dot products and matrix multiplications are performed. They can be improved with more proficient knowledge of the data structures and, specially, of the problem, as there are lots of so called "tricks" that we don't know of that help speed up the optimization and use less resources (for example, sparse matrices manipulation). Lastly, while on our implementations there is no need for the routines to compute the first and second derivatives symbolically, there are not so many in this problem so that it makes a difference.

3.3 Full results

The following is the full Python output of the optimization routine.

```

1 In [8]: print(solution)
2 'xtol' termination condition is satisfied.
3 Number of iterations: 104, function evaluations: 133, CG iterations: 599,
   optimality: 1.83e-07, constraint violation: 0.00e+00, execution time: 0.29 s.
4 barrier_parameter: 2.0480000000000001e-09
5 barrier_tolerance: 2.0480000000000001e-09
6 cg_niter: 599
7 cg_stop_cond: 4
8   constr: [array([1.])]
9   constr_nfev: [133]
10  constr_nhev: [69]
11  constr_njev: [63]
12  constr_penalty: 15407.287523808163
13  constr_violation: 0.0
14  execution_time: 0.2936720848083496
15  fun: 61484.65464046297
16  grad:
17 array([-1.05314317e+04,  1.00060094e+03, -2.22125280e+03,  8.13329381e+02,
18        -1.05852832e+02,  7.31601733e+02,  5.51735226e+03, -1.21838305e+02,
19        -1.40143363e+04,  2.72259723e+03, -4.60324341e+03,  4.48296570e+03,
20        -2.36932573e+03, -8.12422726e+03, -2.14834485e+02, -9.45874490e-11])
21  jac:
22 [array([[0.97647914, -0.09277617,  0.20595557, -0.07541227,  0.00981472,
23         -0.06783445, -0.51157141,  0.0112969 ,  1.29941563, -0.25244045,
24         0.42681482, -0.41566262,  0.21968496,  0.75328205,  0.01991955,
25         0.          ]])]
26  lagrangian_grad:
27 array([-6.64695108e-08,  5.67613370e-08,  4.70404302e-08,  6.35918695e-09,
28         5.68779939e-08, -8.96648089e-10,  3.83806764e-10,  6.66005917e-09,
29         4.45743353e-08, -6.24581844e-08,  4.29736247e-09, -4.31409717e-08,
30        -1.83293196e-07,  1.04246283e-08, -2.83398265e-08, -9.45874490e-11])
31  message: 'xtol' termination condition is satisfied.'
32  method: 'tr_interior_point'
33  nfev: 133
34  nhev: 63
35  niter: 104
36  njev: 63
37  optimality: 1.8329319573240355e-07
38  status: 2
39  tr_radius: 1.3176355152026543e-09
40  v: [array([10785.1067161])]
41  x:
42 array([ 4.88239571e-01, -4.63880870e-02,  1.02977785e-01, -3.77061350e-02,
43        4.90736137e-03, -3.39172227e-02, -2.55785705e-01,  5.64845151e-03,
44        6.49707817e-01, -1.26220227e-01,  2.13407411e-01, -2.07831309e-01,
45        1.09842480e-01,  3.76641023e-01,  9.95977558e-03,  8.95141390e+02])

```

Listing 9: Full Python output.

On a final note, notice that the constrain is not violated: *constraint violation* : 0.00e + 00.

4 Code

Here is the full code, which has also been attached to this document in conjunction with the "Death rate" data set in case it needs to be tried.

```
1 # Authors: Cristina Aguilera, Jesus Antonanzas
2 #          GCED, OM – grup 11.
3 #
4 # Python implementation of ridge – or Tikhonov – regularization:
5 #
6 # min          1/2(A*w + gamma - y)'*(A*w + gamma - y)
7 # subject to   w'*w <= t
8 #
9 # Where 'A' is a set of observations, 'gamma' a vector of real numbers
10 # and 'y' the target variable.
11
12 import numpy as np
13 from scipy.optimize import minimize
14 from scipy.optimize import NonlinearConstraint
15
16 # ===== INPUT DATA =====
17 # Read dataset.
18 # Observations by rows and variables by columns.
19 # Target variable needs to be the last column.
20
21 # A = np.loadtxt('octane_rating.txt', usecols=(1,2,3,4,5))
22 A = np.loadtxt('deathrate_instance_python.dat')
23
24 # Copy target variable.
25 y = A[:, len(A[0]) - 1]
26
27 # Delete target variable from dataset.
28 A = A[:, 0: len(A[0]) - 1]
29
30 # Declare dimensions:
31 #   m = # of observations.
32 #   n = # of variables.
33 m = A.shape[0]
34 n = A.shape[1]
35
36 # ===== PRELIMINARY COMPUTATIONS =====
37
38 # First, declare variables that remain constant
39 # and come in handy during calculations.
40
41 At = A.transpose()
42 AtA = np.matmul(At, A)
43 Aty = np.matmul(At, y)
44 ones_m = np.ones(m)
45 ones_dot_y = np.dot(ones_m, y)
46
47 # Gradient of the restriction:
48 upperleft = np.matmul(At, A)
49 upperright = np.matmul(At, np.ones((m, 1)))
50 tophalf = np.append(upperleft, upperright, axis = 1)
51 bottomhalf = np.append(upperright, m)
52 const_grad_var = np.append(tophalf, [bottomhalf], axis = 0)
53
54 # Hessian of the constraint
55 cons_hess_var = np.diag(np.append(np.full((n), 2), 0))
56
57
58 # ===== EVALUATION FUNCTIONS =====
```

```

59
60 def obj_f(w):
61     # Pos 15 of w contains gamma, all other positions are w_i, i = 0...n-1 (num.
62     # var).
63     gamma = w[n]
64     w = w[0:n]
65     # gamma = gamma * np.ones((m,1)) veure si funciona sense
66     p1 = np.dot(A, w) + gamma * ones_m - y
67     return 1/2 * np.dot(p1, p1)
68
69 def obj_grad(w):
70     # Pos 16 of w contains gamma, all other positions are w_i, i = 0...n-1 (num.
71     # var).
72     gamma = w[n]
73     w = w[0:n]
74     grad = np.zeros(n+1)
75     AtAw = np.matmul(AtA, w)
76     grad[:-1] = AtAw + np.matmul(At, gamma*ones_m) - Aty
77     grad[-1] = np.dot(np.matmul(w, transpose()), At) + m*gamma - ones_dot_y
78     return grad
79
80 def obj_hess(w):
81     # Returned value remains constant, so just calculate once at the initialisation
82     # .
83     return const_grad_var
84
85 def cons(w):
86     w = w[0:n]
87     return np.dot(w, transpose()), w
88
89 def cons_grad(w):
90     return np.append(2*w[0:n], 0)
91
92 def cons_hess(w, v):
93     # Returns sum of all constraint Hessians times their respective Lagrange
94     # multipliers.
95     return v[0]*cons_hess_var
96
97 # ===== MINIMIZATION =====
98
99 # Declare initial values.
100 w0 = np.zeros(n+1)
101 t = 1
102
103 # Declare constraint.
104 nonlinear_constraint = NonlinearConstraint(cons, lb = -np.inf, ub = t,
105                                           jac=cons_grad, hess=cons_hess)
106
107 # Variables are not bounded. We could write
108 # bounds = scipy.optimize.Bounds(-np.inf, np.inf)
109 # but increases computation time.
110
111 # Minimization.
112 solution = minimize(obj_f, w0, method='trust-constr', jac=obj_grad,
113                    hess=obj_hess, constraints=[nonlinear_constraint],
114                    options={'verbose': 1})
115
116 # Display the solution.
117 print(solution)

```

Listing 10: Full Python code.