

SDS 385 Exercise Set 01

Kevin Song

September 3, 2016

1 WLS Objective Function

The weighted least squares problem is described by the equation

$$\hat{\beta} = \arg \min_{\beta \in \mathbb{R}^P} \sum_{i=1}^N \frac{w_i}{2} (y_i - x_i^T \beta)^2,$$

where w_i is the weight of the i -th observation, y_i is the i -th response, and x_i is the i -th row of X , an $N \times P$ matrix.

We can rewrite the objective function in a matrix-vector form in a few steps. First, we note that this is a dot product between the vectors w and $y - X\beta$:

$$\sum_{i=1}^N \frac{w_i}{2} (y_i - x_i^T \beta)^2 = w^T (y - X\beta)$$

Unfortunately, this isn't quite valid, since the exponent denotes an element-wise square, which is something that cannot be done with standard vector operations.

Instead, we convert to quadratic form:

$$\frac{1}{2} w^T (y - X\beta)^2 = \frac{1}{2} (y - X\beta)^T W (y - X\beta)$$

where W is the diagonal matrix with the entries of w along its diagonal.

While I haven't made a formal proof of correctness of this transformation, I offer the following validity argument: define $p = y - X\beta$. If we say that the diagonal elements of W are the elements w , doing the second multiplication first yields the vector q , where $q_i = w_i p_i$. Then we get that the end result is $\sum_i p_i w_i p_i$, which is what we wanted.

Expanding this product, we find that

$$\frac{1}{2} (y - X\beta)^T W (y - X\beta) = \frac{1}{2} (y^T W y - y^T W X \beta - \beta^T X^T W y + \beta^T X^T W X \beta)$$

In this expansion, we find that the first term is useless to us: since we change the value the objective function takes by altering β , and β is not in this term, it does not affect the optimization. We throw it out.

The second and third terms are transposes of each other. As luck would have it, they are also scalar terms. Since $a = a^T$ when a is a scalar, we combine the two terms.

Finally, we recognize that the final term is the quadratic form of β with the matrix $X^T W X$.

This gives us the following objective function:

$$f(\beta) = (y^T W X) \beta + \frac{1}{2} \beta^T X^T W X \beta$$

Note that this is a quadratic form, $\frac{1}{2} x^T A x + b^T x + c$, with $A = X^T W X$, $b = y^T W X$, and $c = 0$.

We know that quadratic forms are minimized by $Ax - b = 0$, or $Ax = b$. Thus, we conclude that $\hat{\beta}$ can be found by solving

$$(X^T W X) \hat{\beta} = y^T W X$$

or equivalently, by

$$(X^T W X) \hat{\beta} = X^T W y$$

2 Solutions to the minimization

The inversion method is a potentially awful way to solve the system. If the entries of X are all relatively similar to each other and the confidence values differ wildly, the computed inverse may not be the inverse at all! More generally, if the condition number of $X^T W X$ is large, then $X^{-1} X$ might not be equal to the identity matrix!

Since we cannot always use matrix inversion, we can instead use a matrix factorization method. In this case, we choose an LU method because I'm not sure if Choleky's conditions will hold here:

```
1  [L,U] = LU_factor(X^T W X)
   z = X^T W y
3  c = triangle_solve(Lc=z)
   x = triangle_solve(Ux=b)
5  beta = x
```

3 Code

Code for solving with inversion:

```
1  ## This solves (X^T W X) B = X^T W y by using matrix inversion on the LHS
   ## Notation: I treat this as an Ax = b problem, so
3  ##   A = X^T W X
   ##   b = X^T W y
5  ##   x = beta

7  import numpy as np
   import scipy as sp
9  import random

11 def solve(A,b):
13     solution = A.I * b
   return solution
```

inversionsolve.py

Code for solving with LU factorization:

```

2  ## This solves  $(X^T W X) B = X^T W y$  by using matrix factorization on the LHS
3  ## Notation: I treat this as an  $Ax = b$  problem, so
4  ##   A =  $X^T W X$ 
5  ##   b =  $X^T W y$ 
6  ##   x = beta
7
8  import numpy as np
9  import scipy as sp
10 from scipy import linalg
11 import random
12
13 def solve(A,b):
14     LUFact = linalg.lu_factor(A)
15     solution = linalg.lu_solve(LUFact, b)
16
17     return solution

```

factorsolve.py

Unfortunately, due to the limitations of my current system (a celeron with 8GB of memory), I was not able to do particularly large tests, since even relatively small timing tests took up a significant amount of time.

The results for the three tests I did manage to finish are shown below. These were done using iPython3's `%timeit` magic. The elements of the matrices are generated random uniform in $[0, 50)$.

Rows	Columns	Iterations	Factorization	Inversion
400	20	1000	2.44	2.54
2000	40	100	18.2	19.7
2000	200	100	33.3	36.3

Table 1: Timings for factorization and inverse-based solutions in python. The times are for a single solve, averaged over the number of iterations (e.g. 5s with 100 iterations means the run took 500 seconds total). All times are in seconds.

4 Sparse Matrices

Converting the code to use sparse matrix representation is pretty simple to do with python. Simply use the sparse solver provided by `scipy.sparse.linalg` to solve the problem. Stacking this solver up against the previous solver, we get the following values:

Rows	Columns	Sparsity	Sparse Solver (ms)	Inversion (ms)	Factorization (ms)
4000	200	0.05	1.63	1.15	0.58
4000	400	0.05	13.3	4.33	2.05
4000	800	0.05	65.1	23.4	10.9
4000	1600	0.05	383	151	52.4
4000	200	0.15	1.77	1.15	0.58
4000	400	0.15	14	4.17	2.02
4000	800	0.15	68.4	23.5	10.2
4000	1600	0.15	394	157	50.9

Table 2: Timings for factorization and inverse-based, as well as sparse solver solutions in python. The non-sparse solvers were run by first converting the sparse matrix representation to a dense one, then using the routines provided by the previous test. All tests were done over 100 iterations.

5 MLE

The MLE form is:

$$\begin{aligned}
l(\beta) &= -\log \left\{ \prod_{i=1}^N p(y_i; \beta) \right\} \\
&= -\sum_{i=1}^N \log p(y_i; \beta) \\
&= -\sum_{i=1}^N \log \left(\binom{m_i}{y_i} (w_i)^{y_i} (1 - w_i)^{m_i - y_i} \right) \\
&= -\sum_{i=1}^N \log \binom{m_i}{y_i} - \sum_{i=1}^N \log w_i^{y_i} - \sum_{i=1}^N \log (1 - w_i)^{m_i - y_i} \\
&= -\sum_{i=1}^N \log \binom{m_i}{y_i} - y_i \sum_{i=1}^N \log w_i - (m_i - y_i) \sum_{i=1}^N \log (1 - w_i)
\end{aligned}$$

Again, since the binomial coefficient is not important in the context of optimizing this function over β , we ignore it.

This gives us the final equation

$$-l(\beta) = y_i \sum_{i=1}^N \log w_i + (m_i - y_i) \sum_{i=1}^N \log (1 - w_i)$$

If we define $l_i(\beta)$ as

$$-l_i(\beta) = \sum_{i=1}^N \log w_i - \sum_{i=1}^N \log (1 - w_i)$$

we can express $l(\beta) = \sum l_i(\beta)$. Then to find the gradient of l , we need only find the gradient of each l_i .

$$-\nabla l_i(\beta) = \nabla_\beta \log w_i + \nabla_\beta \log (1 - w_i)$$

We note that $\nabla_\beta l_i = (w_i)(1 - w_i)x_i$. This allows us to apply the chain rule to the above form, yielding

$$\begin{aligned}
-\nabla l_i(\beta) &= y_i \nabla_\beta \log w_i + (m_i - y_i) \nabla_\beta \log (1 - w_i) \\
&= y_i \frac{1}{w_i} \nabla_\beta w_i + (m_i - y_i) \frac{1}{1 - w_i} \nabla_\beta (1 - w_i) \\
&= y_i \frac{w_i(1 - w_i)}{w_i} x_i - (m_i - y_i) \frac{w_i(1 - w_i)}{1 - w_i} x_i \\
&= y_i(1 - w_i)x_i - (m_i - y_i)w_i x_i \\
&= (y_i - y_i w_i - m_i w_i + y_i w_i)x_i \\
&= (y_i - m_i w_i)x_i
\end{aligned}$$

Finally, this gives us that the gradient of the entire MLE function is

$$\nabla(\beta) = - \sum_{i=1}^N (y_i - m_i w_i) x_i$$

6 Steepest Descent Code

This code has one major quirk: instead of evaluating the gradient and log-likelihood at β_i in a single function, a generator is used to produce a gradient/likelihood function beforehand. This function is then used to compute the desired quantities.

```

1 import math
2 import numpy as np
3 import csv
4 import sys
5 from numpy import linalg as LA
6 import pdb
7
8 bump = 0.00000001 # A tiny bump for some values that really should not be zero
9 smallest_safe_exponent = math.log(sys.float_info.min) + 3
10 largest_safe_exponent = math.log(sys.float_info.max) - 3
11
12 def safe_exp(val):
13     """Calculates the "safe exponential" of a value. If the computed exponential would
14     be too large, it replaces it with a safe value."""
15
16     if val > largest_safe_exponent:
17         print("[WARN]: Exponential term was capped to avoid overflow. This may cause
18         divergence.")
19         return math.exp(largest_safe_exponent)
20     elif val < smallest_safe_exponent:
21         print("[WARN]: Exponential term was capped to avoid underflow. This may cause
22         divergence.")
23         return math.exp(smallest_safe_exponent)
24     else:
25         return math.exp(val)
26
27 def calc_likelihood_function(X,y,m):
28     """Gives a function to determine the likelihood in the inverse logit method.
29     Returns a function which takes in beta and returns the likelihood as a float."""
30
31     def likelihood(B):
32         result = 0
33
34         xvecs = [ xs for xs in X ] # Length Samples
35         exponents = [ np.dot(x,B) for x in xvecs ]
36         exptersms = [ safe_exp(z) for z in exponents ]
37         weights = [ 1.0 / (1.0 + e) for e in exptersms ]
38
39         for i in range(len(xvecs)):
40             result = np.asscalar(y[i]) * math.log(weights[i])
41             result = np.asscalar(m[i] - y[i]) * math.log(1 - weights[i])
42
43         return result
44     return likelihood
45
46 def zoom(calc_deriv, calc_obj, a, c):
47     #Use naive bisection to find trial step
48     (alo,ahi) = a
49     (c1,c2) = c
50
51     i = 1
52     while True:

```

```

51     aj = (ahi + alo) / 2
52     i += 1
53
54     if calc_obj(aj) > calc_obj(0) + c1 * aj * calc_deriv(0) or calc_deriv(aj) >
calc_deriv(alo):
55         ahi = aj
56     else:
57         if abs(calc_deriv(aj)) <= c2 * calc_deriv(0):
58             return aj
59         if calc_deriv(aj) * (ahi - alo) >= 0:
60             ahi = alo
61         alo = aj
62
63     if i > 30:
64         return aj
65
66 def line_search(gradFunc, objFunc, guess, amax, c):
67     (c1, c2) = c
68     i = 1
69     found = False
70
71     searchDir = gradFunc(guess)
72
73     # Helper function to calculate scalar derivatives
74     def calc_deriv(scal):
75         g = gradFunc(guess + searchDir * scal)
76         return np.dot(g.T, searchDir)
77
78     # Helper function to calculate objective fn values
79     def calc_obj(scal):
80         return objFunc(guess + scal * searchDir)
81
82     a = 0.01 * amax
83     a_last = 0
84     a_min = a # The smallest a value so far
85
86     objZero = calc_obj(0)
87     derivZero = calc_deriv(0)
88
89     while True:
90         print(i, a)
91         if calc_obj(a) > objZero + c1 * a * derivZero or (i > 1 and calc_obj(a) >= calc_obj(
a_last)):
92             astar = zoom(calc_deriv, calc_obj, (a_last, a), c)
93             return astar
94         if abs(calc_deriv(a)) <= c2 * derivZero:
95             return a
96         if calc_deriv(a) > 0:
97             astar = zoom(calc_deriv, calc_obj, (a, a_last), c)
98             return astar
99         a_last = a
100         a = a + (amax - a) * 0.1
101
102     i += 1
103
104     if i > 15 or a - a_last < 1000 * sys.float_info.min:
105         return a
106
107 def calc_grad_function(X, y, m):
108     """Calculates the gradient of the inverse logit MLE given the parameters.
109     Return is a lambda function which takes a single parameter and returns float."""
110
111     # X is a feature matrix: columns are features, rows are entries. Dim: samples x features
112     # y is a response vector: one column of responses Dim: samples x 1
113     # m is a trials vector Dim: samples x 1
114
115     def grad(B):
116         # B should be a col vector with length = # features

```

```

117         xvecs = [ xs for xs in X ]                               # Length Samples
119         exponents = [ np.dot(x,B) for x in xvecs ]
121         expters = [ safe_exp(z) for z in exponents ]
123         weights = [ 1.0 / (1.0 + e) for e in expters ]

125         W = np.matrix(np.diag(np.array(weights)))

127         gradient = X.T * (y - W * m)
129         return gradient

131     def solve(params, initial_guess, converge_step):
133         (X,y,m) = params
135         # A function which calculates the gradient at a point
137         grad_op = calc_grad_function(X,y,m)
139         # A function which calculates the likelihood at a point
141         llh_op = calc_likelihood_function(X,y,m)

143         delta = sys.float_info.max
145         guess = initial_guess

147         # For storing likelihoods (for tracking convergence)
149         likelihood_record = []

151         ## Main Steepest Descent Loop
153         while delta > converge_step:
155             oldGuess = guess

157             grad = grad_op(guess)
159             step = line_search(grad_op, llh_op, guess, 0.5, (0.0001, 0.9))
161             print(step)

163             guess = guess - grad * step

165             delta = abs(llh_op(oldGuess) - llh_op(guess))

167             likelihood_record.append(delta)

169             print(delta)

171         return (guess, likelihood_record)

```

steepestdescent.py

7 Quadratic Forms Again

The second-order Taylor approximation of $l(\beta)$ about β_0 is given by Taylor's Theorem:

$$\begin{aligned}
 q(\beta; \beta_0) &\approx l(\beta_0) + \nabla l(\beta_0)^T (\beta - \beta_0) + \frac{1}{2} (\beta - \beta_0)^T \nabla^2 l(\beta_0) (\beta - \beta_0) \\
 &= l(\beta_0) + \nabla l(\beta_0)^T (\beta - \beta_0) + \frac{1}{2} (\beta - \beta_0)^T H_l(\beta_0) (\beta - \beta_0) \\
 &= \left(l(\beta_0) - \nabla l(\beta_0)^T \beta_0 + \frac{1}{2} \beta_0^T H_l(\beta_0) \beta_0 \right) \\
 &\quad + \nabla l(\beta_0)^T \beta - \beta_0^T H_l(\beta_0) \beta + \frac{1}{2} \beta^T H_l(\beta_0) \beta
 \end{aligned}$$

We can get this into the desired quadratic form by completing the square. To derive the form for this, we write the form we have and the form we would like, and attempt to get the two to agree. Assuming that M and C are symmetric,

$$\begin{aligned} a + b^T x + \frac{1}{2} x^T C x &= \frac{1}{2} (x - y)^T M (x - y) \\ &= \frac{1}{2} (x^T M x - x^T M y - y^T M x + y^T M y) + v \\ \frac{1}{2} x^T C x + b^T x + a &= \left(\frac{1}{2} x^T M x \right) - y^T M x + \frac{1}{2} (y^T M y + v) \end{aligned}$$

Comparing the terms on the left and right, this suggests that

$$\begin{aligned} M &= C \\ -y^T M &= b^T \implies y = -M^{-1} b = -C^{-1} b \\ v &= a - \frac{1}{2} y^T M y = a - \frac{1}{2} b^T C^{-1} b \end{aligned}$$

Applying this to our problem, we get that

$$\begin{aligned} M &= H_l(\beta_0) \\ y &= H_l(\beta_0)^{-1} \nabla l(\beta_0) \\ v &= l(\beta_0) - \frac{1}{2} \nabla l(\beta_0)^T H_l(\beta_0)^{-1} \nabla l(\beta_0) \end{aligned}$$

which allows us to write

$$q(\beta; \beta_0) = \frac{1}{2} ()$$