

SDS 385 Exercise Set 01

Kevin Song

August 30, 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
2 ## Notation: I treat this as an Ax = b problem, so
3 ## A = X^T W X
4 ## b = X^T W y
5 ## x = beta
6
7 import numpy as np
8 import scipy as sp
9 import random
10
11 def solve(A,b):
12
13     solution = A.I * b
14     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:

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

The $\log 1$ term is identically equal to zero, and when we take the gradient, we find that the binomial coefficient also goes to zero (since it does not depend on β .)

If we define $l_i(\beta)$ as

$$l_i(\beta) = y_i \sum_{i=1}^N \log(1 + \exp(x_i^T \beta)) - (m_i - y_i) \log \left(\frac{\exp(x_i^T \beta)}{1 + \exp(x_i^T \beta)} \right) - \log \binom{m_i}{y_i}$$

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 .

$$\begin{aligned}
\nabla l_i(\beta) &= \nabla \log(1 + \exp\{x_i^T \beta\}) - \nabla \log \left(\frac{\exp(x_i^T \beta)}{1 + \exp(x_i^T \beta)} \right) \\
&= y_i x_i \left(\frac{\exp(x_i^T \beta)}{1 + \exp(x_i^T \beta)} \right) - (m_i - y_i) x_i \left(\frac{\exp(x_i^T \beta)}{(\exp(x_i^T \beta) + 1)^2} \right)
\end{aligned}$$

which, using the abbreviation mentioned in the assignment, is equivalent to

$$y_i w_i \exp(x_i^T \beta) x_i - (m_i - y_i) w_i^2 \exp(x_i^T \beta) x_i$$

I have no idea if this is even close to correct.