

CS189 HW4

David Chen

February 2022

Contents

1	Write-up and Honor Code	2
2	Logistic Regression with Newton's Method	3
2.1	Part 1	3
2.2	Part 2	4
2.3	Part 3	5
2.4	Part 4	6
2.4.1	Section (a)	6
2.4.2	Section (b)	6
2.4.3	Section (c)	7
2.4.4	Section (d)	7
3	Wine Classification with Logistic Regression	8
3.1	Part 1	8
3.2	Part 2	9
3.3	Part 3	10
3.4	Part 4	11
3.5	Part 5	12
3.6	Part 6	13
4	A Bayesian Interpretation of Lasso	14
4.1	Part 1	14
4.2	Part 2	15
5	l_1-regularization, l_2-regularization, and Sparsity	16
5.1	Part 1	16
5.2	Part 2	17
5.3	Part 3	18
5.4	Part 4	19
5.5	Part 5	20

1 Write-up and Honor Code

Collaborated with: N/A.

I certify that all solutions are entirely in my own words and that I have not looked at another student's solutions. I have given credit to all external sources I consulted.

Signed: David Chen

2 Logistic Regression with Newton's Method

2.1 Part 1

First note that the derivative of the logistic function is given by

$$s'(y) = s(y)(1 - s(y)).$$

Therefore, we can use this result to get

$$\begin{aligned}\nabla_{\mathbf{w}} J(\mathbf{w}) &= \nabla_{\mathbf{w}} - \mathbf{y} \cdot \ln s - (1 - \mathbf{y}) \cdot \ln (1 - \mathbf{y}) \\&= \nabla_{\mathbf{w}} - \sum y_i \ln s_i + (1 - y_i) \ln (1 - s_i) \\&= - \sum (\nabla_{\mathbf{w}} \ln s_i) y_i + (\nabla_{\mathbf{w}} \ln (1 - s_i)) (1 - y_i) \\&= - \sum \frac{1}{s_i} (\nabla_{\mathbf{w}} s_i) y_i - \frac{1}{1 - s_i} (\nabla_{\mathbf{w}} s_i) (1 - y_i) \\&= - \sum \frac{1}{s_i} (\nabla_{\mathbf{w}} s(\mathbf{x}_i \cdot \mathbf{w})) y_i - \frac{1}{1 - s_i} (\nabla_{\mathbf{w}} s(\mathbf{x}_i \cdot \mathbf{w})) (1 - y_i) \\&= - \sum \frac{\mathbf{x}_i}{s_i} (s'_i) y_i - \frac{\mathbf{x}_i}{1 - s_i} (s'_i) (1 - y_i) \\&= - \sum \mathbf{x}_i \left(\frac{1}{s_i} s_i (1 - s_i) y_i - \frac{1}{1 - s_i} s_i (1 - s_i) (1 - y_i) \right) \\&= - \sum \mathbf{x}_i ((1 - s_i) y_i - s_i (1 - y_i)) \\&= - \sum \mathbf{x}_i (y_i - s_i y_i - s_i + s_i y_i) \\&= - \sum \mathbf{x}_i (y_i - s_i) \\&= -X^T (\mathbf{y} - \mathbf{s}).\end{aligned}$$

2.2 Part 2

$$\begin{aligned}
\nabla_{\mathbf{w}}^2 J(\mathbf{w}) &= \nabla_{\mathbf{w}} (\nabla_{\mathbf{w}} J(\mathbf{w})) \\
&= \nabla_{\mathbf{w}} (-X^T (\mathbf{y} - \mathbf{s})) \\
&= \nabla_{\mathbf{w}} X^T \mathbf{s} \\
&= \sum \nabla_{\mathbf{w}} X_i^T s_i \\
&= \sum (\nabla_{\mathbf{w}} s_i) X_i^T \\
&= \sum s_i (1 - s_i) X_i X_i^T \\
&= X^T \Omega X
\end{aligned}$$

where diagonal matrix Ω is given by

$$\Omega = \begin{bmatrix} s_1(1-s_1) & 0 & \dots & 0 \\ 0 & s_2(1-s_2) & \dots & 0 \\ \vdots & \dots & \ddots & \vdots \\ 0 & 0 & \dots & s_n(1-s_n) \end{bmatrix}.$$

2.3 Part 3

Recall that Newton's method updates using vector \mathbf{e} , which is the solution to the linear system $(\nabla_{\mathbf{w}}^2 J(\mathbf{w})) \mathbf{e} = -\nabla_{\mathbf{w}} J(\mathbf{w})$. Therefore, for logistic regression, this becomes

$$\begin{aligned}\mathbf{w} &\leftarrow \mathbf{w} + \mathbf{e} \\ &\leftarrow \mathbf{w} + (X^T \Omega X)^{-1} X^T (\mathbf{y} - \mathbf{s}).\end{aligned}$$

2.4 Part 4

2.4.1 Section (a)

Design matrix X is given by

$$X = \begin{bmatrix} 0.2 & 3.1 & 1 \\ 1.0 & 3.0 & 1 \\ -0.2 & 1.2 & 1 \\ 1.0 & 1.1 & 1 \end{bmatrix}.$$

Therefore,

$$\begin{aligned} X\mathbf{w}^{(0)} &= \begin{bmatrix} 0.2 & 3.1 & 1 \\ 1.0 & 3.0 & 1 \\ -0.2 & 1.2 & 1 \\ 1.0 & 1.1 & 1 \end{bmatrix} \begin{bmatrix} -1 \\ 1 \\ 0 \end{bmatrix} \\ &= \begin{bmatrix} 2.9 \\ 2.0 \\ 1.4 \\ 0.1 \end{bmatrix}. \end{aligned}$$

$\mathbf{s}^{(0)}$ is equal to the logistic function applied element-wise to $X\mathbf{w}^{(0)}$, which gives

$$\begin{aligned} \mathbf{s}^{(0)} &= \begin{bmatrix} s(2.9) \\ s(2.0) \\ s(1.4) \\ s(0.1) \end{bmatrix} \\ &= \begin{bmatrix} 0.9478 \\ 0.8808 \\ 0.8022 \\ 0.5250 \end{bmatrix}. \end{aligned}$$

2.4.2 Section (b)

We first find \mathbf{e} with

$$\begin{aligned} \mathbf{e} &= (X^T \Omega X)^{-1} X^T (\mathbf{y} - \mathbf{s}) \\ &= \begin{bmatrix} 0.0953 \\ 0.5623 \\ -1.6783 \end{bmatrix}. \end{aligned}$$

Following the update rule, we get

$$\begin{aligned}
\mathbf{w}^{(1)} &= \mathbf{w}^{(0)} + \mathbf{e} \\
&= [-0.9047 \quad 1.5623 \quad -1.6783]^T.
\end{aligned}$$

2.4.3 Section (c)

Following the same procedure as above gives

$$\begin{aligned}
X\mathbf{w}^{(1)} &= \begin{bmatrix} 2.9839 \\ 2.1039 \\ 0.3774 \\ -0.8645 \end{bmatrix} \\
\therefore \mathbf{s}^{(1)} &= \begin{bmatrix} s(2.9839) \\ s(2.1039) \\ s(0.3774) \\ s(-0.8645) \end{bmatrix} \\
&= \begin{bmatrix} 0.9518 \\ 0.8913 \\ 0.5932 \\ 0.2964 \end{bmatrix}.
\end{aligned}$$

2.4.4 Section (d)

Following the same procedure as above gives

$$\begin{aligned}
\mathbf{e} &= (X^T \Omega X)^{-1} X^T (\mathbf{y} - \mathbf{s}) \\
&= \begin{bmatrix} 0.1614 \\ 0.3733 \\ -1.1480 \end{bmatrix} \\
\therefore \mathbf{w}^{(2)} &= \mathbf{w}^{(1)} + \mathbf{e} \\
&= [-0.7433 \quad 1.9356 \quad -2.8263]^T.
\end{aligned}$$

3 Wine Classification with Logistic Regression

3.1 Part 1

The only difference between logistic regression and logistic regression l_2 regularization is the added $\lambda\|\mathbf{w}\|^2$ term to the cost function. This gives

$$\begin{aligned} J(\mathbf{w}) &= -\mathbf{y} \cdot \ln \mathbf{s} - (\mathbf{1} - \mathbf{y}) \cdot \ln (\mathbf{1} - \mathbf{s}) + \lambda\|\mathbf{w}\|^2 \\ \therefore \nabla_{\mathbf{w}} J(\mathbf{w}) &= \nabla_{\mathbf{w}} (-\mathbf{y} \cdot \ln \mathbf{s} - (\mathbf{1} - \mathbf{y}) \cdot \ln (\mathbf{1} - \mathbf{s})) + \nabla_{\mathbf{w}} (\lambda\|\mathbf{w}\|^2) \\ &= -X^T (\mathbf{y} - \mathbf{s}) + 2\lambda\mathbf{w} \\ \therefore \mathbf{w} &\leftarrow \mathbf{w} - \varepsilon \nabla_{\mathbf{w}} J(\mathbf{w}) \\ &\leftarrow \mathbf{w} - \varepsilon (-X^T (\mathbf{y} - \mathbf{s}) + 2\lambda\mathbf{w}) \\ &\leftarrow \mathbf{w} + \varepsilon (X^T (\mathbf{y} - \mathbf{s}) - 2\lambda\mathbf{w}). \end{aligned}$$

3.2 Part 2

Included in attached Python Notebook.

3.3 Part 3

Since stochastic gradient descent works with one training example rather than all, it uses a similar update

$$\begin{aligned}w &\leftarrow w - \varepsilon \left(-\mathbf{X}_i^T (y_i - s(\mathbf{X}_i \cdot \mathbf{w})) + 2\lambda \mathbf{w} \right) \\&\leftarrow w + \varepsilon \left(\mathbf{X}_i^T (y_i - s(\mathbf{X}_i \cdot \mathbf{w})) - 2\lambda \mathbf{w} \right).\end{aligned}$$

3.4 Part 4

Included in attached Python Notebook.

3.5 Part 5

Included in attached Python Notebook.

3.6 Part 6

Included in attached Python Notebook.

4 A Bayesian Interpretation of Lasso

4.1 Part 1

By Bayes' Theorem,

$$f\left(\mathbf{w} | (\mathbf{x}_i, y_i)_{i \in [n]}\right) = \frac{f(y_i | \mathbf{x}_i, \mathbf{w}) \cdot f(\mathbf{w} | \mathbf{x}_i)}{f(y_i)}.$$

Since \mathbf{w} is a random parameter, $f(\mathbf{w} | \mathbf{x}_i) = f(\mathbf{w})$. Therefore,

$$f\left(\mathbf{w} | (\mathbf{x}_i, y_i)_{i \in [n]}\right) = \frac{f(y_i | \mathbf{x}_i, \mathbf{w}) \cdot f(\mathbf{w})}{f(y_i)},$$

where $f(y_i)$ is the maximum likelihood estimate of y_i attained by counting the amount of y_i 's and dividing by the total length of y .

4.2 Part 2

$$\begin{aligned}
\max_{\mathbf{w}} l(\mathbf{w}) &= \max_{\mathbf{w}} \ln f(\mathbf{w} | (\mathbf{x}_i, y_i)_{i \in [n]}) \\
&= \max_{\mathbf{w}} \ln \frac{f(y_i | \mathbf{x}_i, \mathbf{w}) \cdot f(\mathbf{w})}{f(y_i)} \\
&= \max_{\mathbf{w}} \ln((y_i | \mathbf{x}_i, \mathbf{w})) + \ln(f(\mathbf{w})) - \ln(f(y_i)) \\
&= \max_{\mathbf{w}} \ln((y_i | \mathbf{x}_i, \mathbf{w})) + \ln(f(\mathbf{w})) \\
&= \max_{\mathbf{w}} \ln \left(\frac{1}{\sigma \sqrt{2\pi}} e^{-\frac{(y - \mathbf{w} \cdot \mathbf{x})^2}{2\sigma^2}} \right) + \ln \left(\frac{1}{2b} e^{-\frac{\|\mathbf{w}\|_1}{b}} \right) \\
&= \max_{\mathbf{w}} \ln \left(\frac{1}{\sigma \sqrt{2\pi}} \right) \left(-\frac{(y - \mathbf{w} \cdot \mathbf{x})^2}{2\sigma^2} \right) + \ln \left(\frac{1}{2b} \right) \left(-\frac{\|\mathbf{w}\|_1}{b} \right) \\
&= \max_{\mathbf{w}} -(y - \mathbf{w} \cdot \mathbf{x})^2 - \lambda \|\mathbf{w}\|_1 \\
&= \min_{\mathbf{w}} (y - \mathbf{w} \cdot \mathbf{x})^2 + \lambda \|\mathbf{w}\|_1,
\end{aligned}$$

where $\lambda = \frac{\ln\left(\frac{1}{\sigma \sqrt{2\pi}}\right)}{\ln\left(\frac{1}{2b}\right)} \frac{2\sigma^2}{b}$.

5 l_1 -regularization, l_2 -regularization, and Sparsity

5.1 Part 1

Included in attached Python Notebook.

5.2 Part 2

Proof.

$$\begin{aligned}
J_1(\mathbf{w}) &= \|X\mathbf{w} - \mathbf{y}\|^2 + \lambda \|\mathbf{w}\|_1 \\
&= (X\mathbf{w} - \mathbf{y})^T (X\mathbf{w} - \mathbf{y}) + \lambda \sum_i w_i \\
&= (\mathbf{w}^T X^T - \mathbf{y}^T) (X\mathbf{w} - \mathbf{y}) + \lambda \sum_i w_i \\
&= \|\mathbf{y}\|^2 + \|X\mathbf{w}\|^2 - 2\mathbf{y}^T X\mathbf{w} + \lambda \sum_i w_i \\
&= \|\mathbf{y}\|^2 + \sum_i \left(w_i^2 \|\mathbf{x}_{*i}\|^2 - 2\mathbf{y}^T \mathbf{x}_{*i} w_i + \lambda w_i \right) \\
&= \|\mathbf{y}\|^2 + \sum_i \left(w_i^2 \|nI_{*i}\|^2 - 2\mathbf{y}^T \mathbf{x}_{*i} w_i + \lambda w_i \right) \\
&= \|\mathbf{y}\|^2 + \sum_i \left(w_i^2 n^2 - 2\mathbf{y}^T \mathbf{x}_{*i} w_i + \lambda w_i \right) \\
&= \|\mathbf{y}\|^2 + \sum_i f(\mathbf{x}_{*i}, w_i).
\end{aligned}$$

□

5.3 Part 3

$$\begin{aligned}\frac{\partial f(\mathbf{x}_{*i}, \mathbf{w}_i)}{\partial \mathbf{w}_i} &= \frac{\partial}{\partial \mathbf{w}_i} (\mathbf{w}_i^2 n^2 - 2\mathbf{y}^T \mathbf{x}_{*i} \mathbf{w}_i + \lambda \mathbf{w}_i) \\ &= 2\mathbf{w}_i n^2 - 2\mathbf{y}^T \mathbf{x}_{*i} + \lambda \\ \therefore \mathbf{w}_i^* &= \frac{2\mathbf{y}^T \mathbf{x}_{*i} - \lambda}{2n^2}.\end{aligned}$$

This gives that the sign of \mathbf{w}_i^* is equal to the sign of $2\mathbf{y}^T \mathbf{x}_{*i} - \lambda$, as the denominator is always guaranteed to be positive. In addition, since

$$\begin{aligned}\frac{\partial^2 f(\mathbf{x}_{*i}, \mathbf{w}_i)}{\partial \mathbf{w}_i^2} &= \frac{\partial}{\partial \mathbf{x}_{*i}} (2\mathbf{w}_i n^2 - 2\mathbf{y}^T \mathbf{x}_{*i} + \lambda) \\ &= 2n^2 \\ &> 0,\end{aligned}$$

$f(\mathbf{x}_{*i}, \mathbf{w}_i)$ is a convex function and \mathbf{w}_i^* is guaranteed to be a minimum.

5.4 Part 4

Changing from $l1$ -regularization to $l2$ -regularization changes $f(\mathbf{x}_{*i}, \mathbf{w}_i)$ to

$$f(\mathbf{x}_{*i}, \mathbf{w}_i) = \mathbf{w}_i^2 n^2 - 2\mathbf{y}^T \mathbf{x}_{*i} \mathbf{w}_i + \lambda \mathbf{w}_i^2.$$

Therefore, the first derivative becomes

$$\begin{aligned} \frac{\partial f(\mathbf{x}_{*i}, \mathbf{w}_i)}{\partial \mathbf{w}_i} &= \frac{\partial}{\partial \mathbf{w}_i} (\mathbf{w}_i^2 n^2 - 2\mathbf{y}^T \mathbf{x}_{*i} \mathbf{w}_i + \lambda \mathbf{w}_i^2) \\ &= 2\mathbf{w}_i n^2 - 2\mathbf{y}^T \mathbf{x}_{*i} + 2\lambda \mathbf{w}_i \\ \therefore \mathbf{w}_i^\# &= \frac{\mathbf{y}^T \mathbf{x}_{*i}}{n^2 + \lambda}, \end{aligned}$$

which implies that $\mathbf{w}_i^\# = 0$ if $\mathbf{y}^T \mathbf{x} = 0$. Furthermore, this is only guaranteed to hold if $f(\mathbf{x}_{*i}, \mathbf{w}_i)$ is convex, which occurs when

$$\begin{aligned} \frac{\partial^2 f(\mathbf{x}_{*i}, \mathbf{w}_i)}{\partial \mathbf{w}_i^2} &= \frac{\partial}{\partial \mathbf{w}_i} (2\mathbf{w}_i n^2 - 2\mathbf{y}^T \mathbf{x}_{*i} + 2\lambda \mathbf{w}_i) \\ &= 2n^2 + 2\lambda \\ \therefore n^2 &> -\lambda, \end{aligned}$$

which we add as another necessary and sufficient condition.

5.5 Part 5

\mathbf{w}^* is more likely to be sparse, as it is more unlikely for $\mathbf{y}^T \mathbf{x}_{*i}$ to be equal to a specific value λ rather than 0. The exception to this is when $\lambda = 0$, which causes \mathbf{w}^* and $\mathbf{w}^\#$ to have the same sparsity.