

HW3 Report

Section 1

Q1:

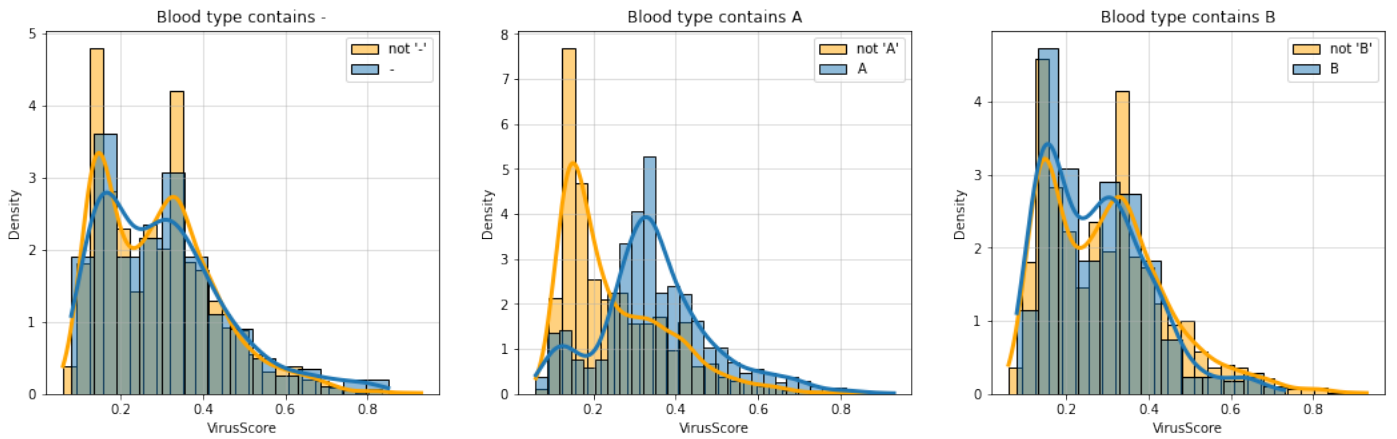


Figure 1: KDE plots of `VirusScore` conditioned different conditions of `blood_type`

Q2:

In figure 1 in the plot of A versus not A, we observe that the groups of patients with and without "A" in their blood types are mostly separable along a boundary that is approximately the `VirusScore` of 0.225.

Therefore, the condition of contains/does not contain A would be most informative for learning `VirusScore`. As it turns out, we decided already in hw1 to create this feature.

Q3:

$$\begin{aligned}
\frac{\partial}{\partial b} \mathcal{L}(\mathbf{w}, b) &\stackrel{a}{=} \frac{\partial}{\partial b} \frac{1}{m} \sum_{i=1}^m (\mathbf{w}^\top \mathbf{x}_i + b - y_i)^2 \stackrel{b}{=} \frac{1}{m} \frac{\partial}{\partial b} \sum_{i=1}^m (\mathbf{w}^\top \mathbf{x}_i + b - y_i)^2 \stackrel{c}{=} \frac{1}{m} \sum_{i=1}^m \frac{\partial}{\partial b} (\mathbf{w}^\top \mathbf{x}_i + b - y_i)^2 \stackrel{d}{=} \\
&\stackrel{d}{=} \frac{1}{m} \sum_{i=1}^m 2 (\mathbf{w}^\top \mathbf{x}_i + b - y_i) \cdot (1) \stackrel{e}{=} \frac{2}{m} \sum_{i=1}^m (\mathbf{w}^\top \mathbf{x}_i + b - y_i) \stackrel{f}{=} \frac{2}{m} \left[mb + \sum_{i=1}^m (\mathbf{w}^\top \mathbf{x}_i - y_i) \right] \stackrel{g}{=} \\
&\stackrel{g}{=} 2b + \frac{2}{m} \sum_{i=1}^m (\mathbf{w}^\top \mathbf{x}_i - y_i) \Rightarrow \frac{\partial}{\partial b} \mathcal{L}(\mathbf{w}, b) = 2b + \frac{2}{m} \sum_{i=1}^m (\mathbf{w}^\top \mathbf{x}_i - y_i)
\end{aligned}$$

a : Definition of $\mathcal{L}(\mathbf{w}, b)$

b : $\frac{1}{m}$ is scalar

c : Derivative of a sum is the sum of derivatives

d : Derivative of $(\mathbf{w}^\top \mathbf{x}_i + b - y_i)^2$ w.r.t b

e : 2 is scalar

f : Sum of b

g : Removing b from the sum

Q4:

Residuals of analytical and numerical gradients

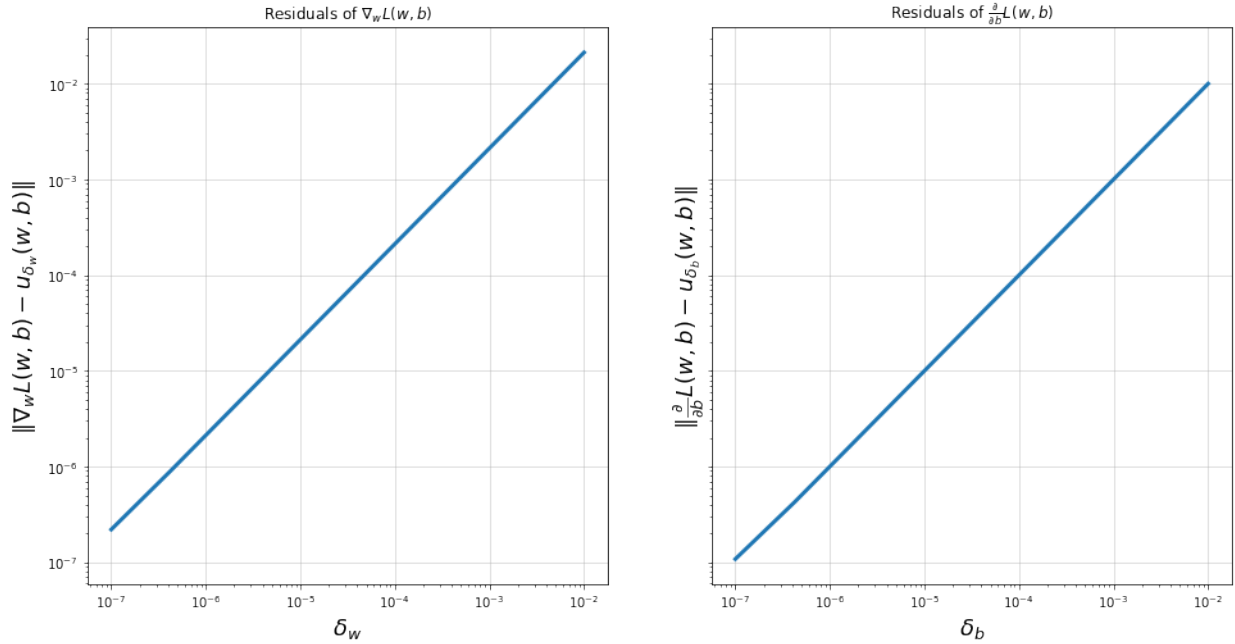


Figure 2: Plot of Residuals of analytical and numerical gradients

As we can see in figure 2, the difference between the analytic and numerical gradients increases in a monotonic fashion as the value of δ increases. This is logical, as δ is the differential size used in the definition of the numerical gradient, and therefore a smaller δ equates to a more precise estimation of the analytic gradient by the numerical gradient.

Q5:

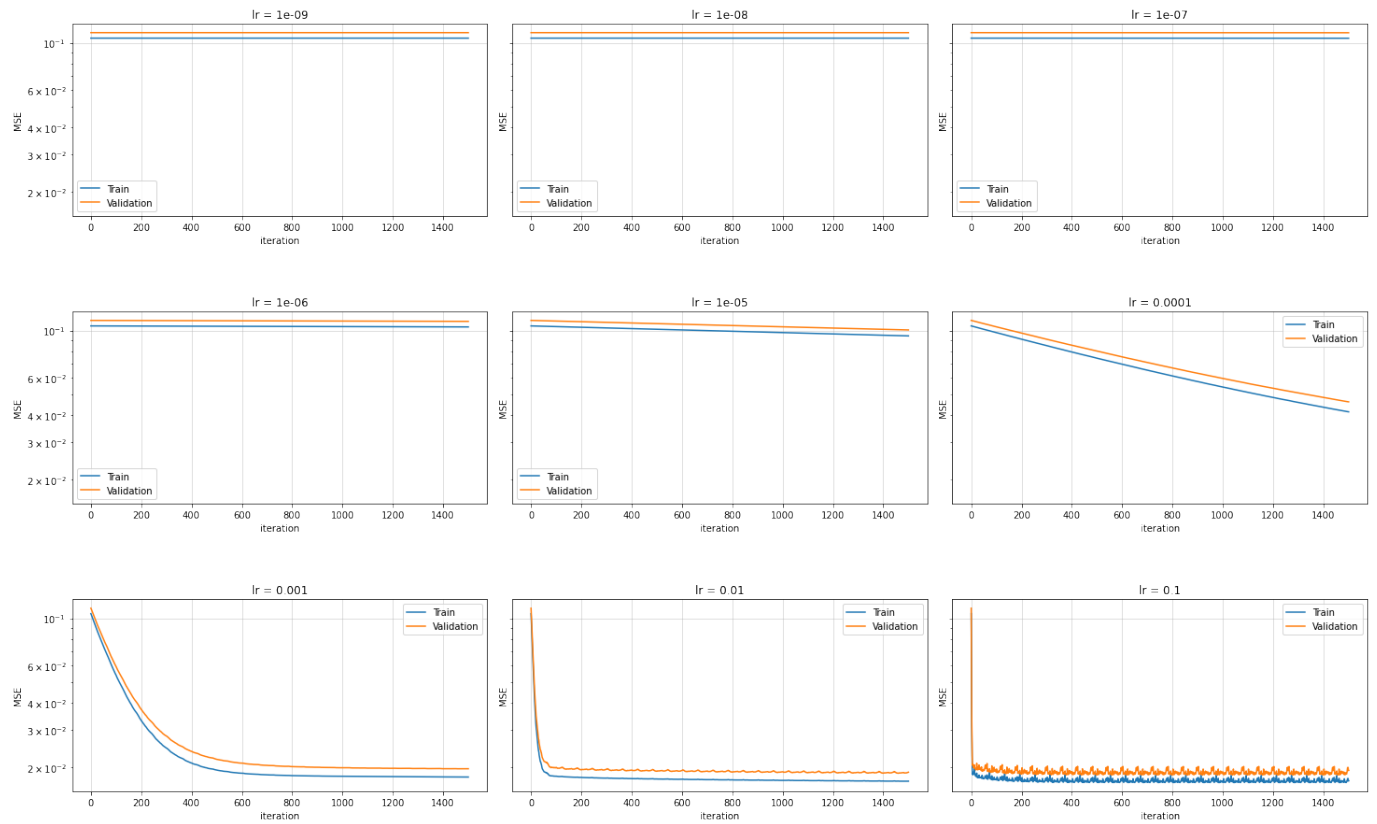


Figure 3: Graphs of Training and Validation Losses as Functions of Iteration Number for Different Learning Rates

We can see in figure 3 that for smaller lr , for those that converge, the convergence is at a higher loss for both the training and validation. This matches the theory, since if the lr is too small, the SGD algorithm is likely to converge to a sub-optimal solution. In addition we observe that for the lr equal to 0.1 the graphs do not converge for both training and validation losses and for lr equal to 0.01 the validation loss does not converge, which fits the theory that says that learning rates that are too high are likely to cause the SGD algorithm to take steps that are too large and thus repeatedly skip-over the optimal solution. This points to 0.001 as being the optimal lr , as both the validation and training losses converge to values that are substantially lower than the next smaller lr .

Q6:

Model	Section	Train MSE	Valid MSE
		Cross Validated	
Dummy	3	0.0204	0.0205

Section 4

Q8:

Model	Section	Train MSE	Valid MSE
		Cross Validated	
Dummy	3	0.0204	0.0205
Ridge linear	4	0.0204	0.0205

Section 6

Q16:

When using a polynomial feature mapping, we can expect the training error to decrease and the validation error to increase. This is because the polynomial mapping will give more flexibility to the regression model to more closely try to fit the training data during training in a polynomial way, thereby leading to lower training error, but also causing overfitting and therefore leading to a higher validation error.