

1 Logistic Regression

Assume that we have n i.i.d. data points $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)$, where each y_i is a binary label in $\{0, 1\}$. We model the posterior probability as a Bernoulli distribution and the probability for each class is the sigmoid function, i.e., $p(y|\mathbf{x}; \mathbf{w}) = q^y(1 - q)^{1-y}$, where $q = s(\mathbf{w}^\top \mathbf{x})$ and $s(\zeta) = \frac{1}{1+e^{-\zeta}}$ is the sigmoid function. Write out the likelihood and log likelihood functions. Comment on whether it is possible to find a closed form maximum likelihood estimate of \mathbf{w} , and describe an alternate approach.

2 Initialization of Weights for Backpropagation

Assume a fully-connected 1-hidden-layer network. Denote the dimensionalities of the input, hidden, and output layers as $d^{(0)}$, $d^{(1)}$, and $d^{(2)}$. That is, the input (which we will denote with a superscript (0)) is a vector of the form $x_1^{(0)}, \dots, x_{d^{(0)}}^{(0)}$. Let g denote the activation function applied at each layer. We will let $S_j^{(l)} = \sum_{i=1}^{d^{(l-1)}} w_{ij}^{(l)} x_i^{(l-1)}$ be the weighted input to node j in layer l , and let $\delta_j^{(l)} = \frac{\partial \ell}{\partial S_j^{(l)}}$ be the partial derivative of the final loss ℓ with respect to $S_j^{(l)}$.

Recall that backpropagation is an efficient method to compute the gradient of the loss function so we can use it for gradient descent. Gradient descent requires the parameters to be initialized to some value(s).

- (a) To better orient yourself with the operations described in this 1-hidden-layer network, draw out a diagram of the layers, including weights, activation functions, and the outputs of each operation during the forward pass. In addition, identify where the partial derivatives $\delta_j^{(l)}$ are calculated during backpropagation.
- (b) Imagine that we initialize every element of each weight $w^{(l)}$ to be the same constant scalar value a . After performing the forward pass, what is the value of $x_j^{(1)}$ in terms of the elements of $\{x_i^{(0)} : i = 1, \dots, d^{(0)}\}$? What is the relationship between each $x_j^{(1)}$?
- (c) Following from the previous part, after the backward pass of backpropagation, compute the values for each member of the set $\{\delta_i^{(1)} : i = 1, \dots, d^{(1)}\}$, assuming we have calculated $\{\delta_j^{(2)} : j = 1, \dots, d^{(2)}\}$. What is the relationship between each $\delta_i^{(1)}$?
- (d) For a reasonable loss function, are all of the $\delta_i^{(2)}$ equal to each other? Why or why not?
- (e) In the previous part, you showed that $w_{ij}^{(2)}$ is different for each j , but for a fixed j , it is the same for each i . In fact, no matter how many subsequent iterations of gradient descent you take, this property will continue to be true. Show why this is the case.
- (f) To solve this problem, we randomly initialize our weights. This is called symmetry breaking. Note that for logistic regression, we don't run into this issue; that is, gradient descent will find the optimal values of the weights even if we initialize them at 0. Explain why this discrepancy exists between our 1-hidden-layer neural network and logistic regression.