

HW3 Written Part

Due 5pm Friday Feb 21, 2020

What is your name?

Problem 1

Suppose I fit a neural network model with the following structure:

- Input layer has 2 features
- One hidden layer with 2 units and a linear activation function
- Output layer has 1 unit and a sigmoid activation function

Show that this model is equivalent to a logistic regression model in the sense that the activation in the output layer could be written as

$$a_1^{(i)[2]} = \frac{\exp(b^* + w_{11}^* x_1^{(i)} + w_{12}^* x_2^{(i)})}{1 + \exp(b^* + w_{11}^* x_1^{(i)} + w_{12}^* x_2^{(i)})}$$

for some parameters b^* , w_{11}^* , and w_{12}^* that are combinations of the biases and weights in all units of the full neural network model. Your answer should give exact formulas for how to calculate b^* , w_{11}^* , and w_{12}^* in terms of the neural network parameters $b_1^{[1]}$, $w_{11}^{[1]}$, $w_{12}^{[1]}$, $b_2^{[1]}$, $w_{21}^{[1]}$, $w_{22}^{[1]}$, $b_1^{[2]}$, $w_{11}^{[2]}$, and $w_{12}^{[2]}$. Comment briefly (1 sentence) on why it is necessary to use non-linear activation functions in neural network models.

Problem 2

Suppose I am working on a classification problem where the response has three classes and I have two input features. I will use a neural network model with the following structure:

- Input layer has 2 features
- One hidden layer has 2 units and a relu activation
- Output layer has 3 units and a softmax activation

My full data set has 100 observations in it.

(a) Give the shapes of each of the following quantities. Use the convention that each observation is in a column of X and each feature is in a row of X . For example, if I am predicting whether an animal is a bird, a cat, or a dog using its weight and its height, the weights and height for the first animal in my data set would be in the first column of X . Also, suppose

we are using a one-hot encoding for the response, so if the first animal in my data set is a dog then I will have $y^{(1)} = \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix}$

- X
- y
- $z^{[1]}$
- $a^{[1]}$
- $b^{[1]}$
- $w^{[1]}$
- $z^{[2]}$
- $a^{[2]}$
- $b^{[2]}$
- $w^{[2]}$
- $\frac{\partial J(b,w)}{\partial a^{[2]}}$
- $\frac{\partial J(b,w)}{\partial z^{[2]}}$
- $\frac{\partial J(b,w)}{\partial b^{[2]}}$
- $\frac{\partial J(b,w)}{\partial w^{[2]}}$
- $\frac{\partial J(b,w)}{\partial a^{[1]}}$
- $\frac{\partial J(b,w)}{\partial z^{[1]}}$
- $\frac{\partial J(b,w)}{\partial b^{[1]}}$
- $\frac{\partial J(b,w)}{\partial w^{[1]}}$

(b) In the backpropagation algorithm, why do we calculate $\frac{\partial J(b,w)}{\partial a^{[2]}}$ before we calculate $\frac{\partial J(b,w)}{\partial z^{[2]}}$? Your answer should involve a formula for how $a^{[2]}$ is calculated and an application of the chain rule.

(c) In the backpropagation algorithm, why do we calculate $\frac{\partial J(b,w)}{\partial z^{[2]}}$ before we calculate $\frac{\partial J(b,w)}{\partial b^{[2]}}$? Your answer should involve a formula for how $z^{[2]}$ is calculated and an application of the chain rule.

Problem 3

Suppose I am working on a classification problem where the response has two classes (say dog and cat) and I have one input feature. In the model statements below, I'm suppressing as much notation as possible.

Our first option for this task is a logistic regression model where $Y^{(i)}$ is encoded as 0 for a dog or 1 for a cat:

$$Y^{(i)} \sim \text{Bernoulli}(a^{(i)})$$
$$a_1^{(i)} = \frac{\exp(b + w_1 x^{(i)})}{1 + \exp(b + w_1 x^{(i)})}$$

However, a reasonable person might also formulate this as a multinomial regression problem using a one-hot encoding of $Y^{(i)} = \begin{bmatrix} 1 \\ 0 \end{bmatrix}$ for a dog or $Y^{(i)} = \begin{bmatrix} 0 \\ 1 \end{bmatrix}$ for a cat:

$$Y^{(i)} \sim \text{Categorical}(a_1^{(i)}, a_2^{(i)})$$
$$a_1^{(i)} = \frac{\exp(b_1 + w_1 x^{(i)})}{\exp(b_1 + w_1 x^{(i)}) + \exp(b_2 + w_2 x^{(i)})}$$
$$a_2^{(i)} = \frac{\exp(b_2 + w_2 x^{(i)})}{\exp(b_1 + w_1 x^{(i)}) + \exp(b_2 + w_2 x^{(i)})}$$

Note that by convention the numbering of classes is 0 and 1 in the logistic regression model, but 1 and 2 in the multinomial regression model. So class 1 in the logistic regression model refers to the same thing as class 2 in the multinomial regression model. This is awkward but I think it'll be more confusing if we change the standard notation. . .

Suppose these models will be estimated by gradient descent, and the parameter values for the two models are initialized so that for any value of x the initial estimated probability of being a cat from the logistic model is equal to the initial estimated probability of being a cat from the multinomial regression model.

(a) Write down the formulas for the updates to b and w for the logistic regression model in terms of $a^{(i)}$, $y^{(i)}$, and $x^{(i)}$, $i = 1, \dots, m$. Note that you don't need to calculate the value of $a^{(i)}$ explicitly in terms of b and w .

(b) Suppose we have two observations with feature, response, and output layer activation values for a logistic regression model as given in the table below. The current parameter values are $b = 1$ and $w = -1$. Find the updated parameter values after one step using a learning rate of $\alpha = 0.1$.

$x^{(i)}$	$y^{(i)}$	$a_1^{(i)}$
1	1	0.5
2	0	0.269

i. Find $\frac{\partial J(b,w)}{\partial z^{[1]}}$. (First, think about what its shape should be.)

ii. Find $\frac{\partial J(b,w)}{\partial b}$. (First, think about what its shape should be.)

iii. Find $\frac{\partial J(b,w)}{\partial w}$. (First, think about what its shape should be.)

iv. Find the updated values of b and w from one gradient descent update step using a learning rate of $\alpha = 0.01$.

(c) Write down the formulas for the updates to b and w for the multinomial regression model in terms of $a^{(i)}$, $y^{(i)}$, and $x^{(i)}$, $i = 1, \dots, m$. Note that you don't need to calculate the value of $a^{(i)}$ explicitly in terms of b and w .

(d) Suppose I have two observations with feature, response, and output layer activation values for a multinomial regression model as given in the table below. My current parameter values are $b = \begin{bmatrix} 0 \\ 1 \end{bmatrix}$ and $w = \begin{bmatrix} 0 \\ -1 \end{bmatrix}$. Find the updated parameter values after one step using a learning rate of $\alpha = 0.1$.

$x^{(i)}$	$y^{(i)}$	$a_1^{(i)}$	$a_2^{(i)}$
1	1	0.5	0.5
2	0	0.731	0.269

i. Find $\frac{\partial J(b,w)}{\partial z^{(1)}}$. (First, think about what its shape should be.)

ii. Find $\frac{\partial J(b,w)}{\partial b}$. (First, think about what its shape should be.)

iii. Find $\frac{\partial J(b,w)}{\partial w}$. (First, think about what its shape should be.)

iv. Find the updated values of b and w from one gradient descent update step using a learning rate of $\alpha = 0.01$.

(e) Based on your answers to parts (a) and (c), argue that if the logistic model and multinomial model currently provide the same probability that each animal is a cat (so that $a_1^{(i)}$ in the logistic regression model is equal to $a_2^{(i)}$ in the multinomial regression model for all observations i), the updates to b and w in the logistic regression model will be the same as the updates to b_2 and w_2 in the multinomial regression model.