

Q1. Logistic Regression

We would like to classify some data. We have N samples, where each sample consists of a feature vector $\mathbf{x} = \{x_1, \dots, x_k\}$ and a label $y = \{0, 1\}$.

We introduce a new type of classifier called logistic regression, which produces predictions as follows:

$$P(Y = 1|X) = h(\mathbf{x}) = s\left(\sum_i w_i x_i\right) = \frac{1}{1 + \exp(-(\sum_i w_i x_i))}$$

$$s(\gamma) = \frac{1}{1 + \exp(-\gamma)}$$

where $s(\gamma)$ is the logistic function, $\exp x = e^x$, and $\mathbf{w} = \{w_1, \dots, w_k\}$ are the learned weights.

Let's find the weights w_j for logistic regression using stochastic gradient descent. We would like to minimize the following loss function for each sample:

$$L = -[y \ln h(\mathbf{x}) + (1 - y) \ln(1 - h(\mathbf{x}))]$$

(a) Find dL/dw_i . Hint: $s'(\gamma) = s(\gamma)(1 - s(\gamma))$.

Use chain rule:

$$\frac{dL}{dw_i} = - \left[\frac{y}{h(\mathbf{x})} s'(\sum_i w_i x_i) x_i - \frac{1 - y}{1 - h(\mathbf{x})} s'(\sum_i w_i x_i) x_i \right]$$

Use hint:

$$\frac{dL}{dw_i} = - \left[\frac{y}{h(\mathbf{x})} h(\mathbf{x})(1 - h(\mathbf{x})) x_i - \frac{1 - y}{1 - h(\mathbf{x})} h(\mathbf{x})(1 - h(\mathbf{x})) x_i \right]$$

Simplify:

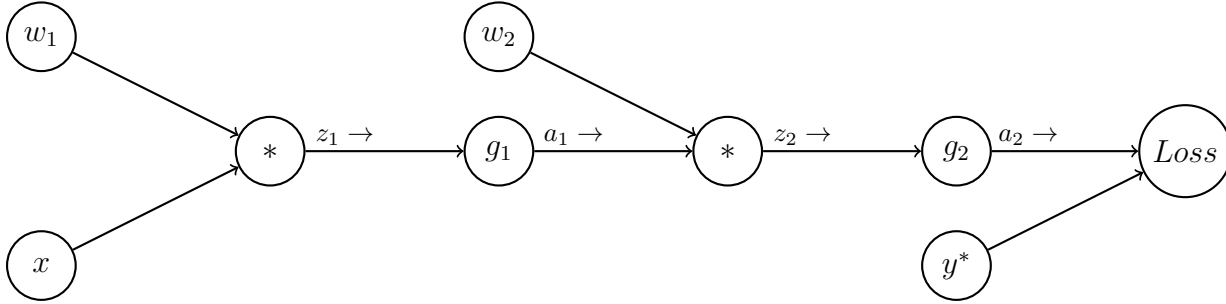
$$\begin{aligned} \frac{dL}{dw_i} &= - [y(1 - h(\mathbf{x})) x_i - (1 - y)h(\mathbf{x}) x_i] \\ &= -x_i [y - yh(\mathbf{x}) - h(\mathbf{x}) + yh(\mathbf{x})] \\ &= -x_i (y - h(\mathbf{x})) \end{aligned}$$

(b) Write the stochastic gradient descent update for w_i . Our step size is η .

$$w_i \leftarrow w_i + \eta x_i (y - h(\mathbf{x}))$$

2 Neural Nets

Consider the following computation graph for a simple neural network for binary classification. Here x is a single real-valued input feature with an associated class y^* (0 or 1). There are two weight parameters w_1 and w_2 , and non-linearity functions g_1 and g_2 (to be defined later, below). The network will output a value a_2 between 0 and 1, representing the probability of being in class 1. We will be using a loss function $Loss$ (to be defined later, below), to compare the prediction a_2 with the true class y^* .



1. Perform the forward pass on this network, writing the output values for each node z_1, a_1, z_2 and a_2 in terms of the node's input values:

$$\begin{aligned} z_1 &= x * w_1 \\ a_1 &= g_1(z_1) \\ z_2 &= a_1 * w_2 \\ a_2 &= g_2(z_2) \end{aligned}$$

2. Compute the loss $Loss(a_2, y^*)$ in terms of the input x , weights w_i , and activation functions g_i :
 Recursively substituting the values computed above, we have:

$$Loss(a_2, y^*) = Loss(g_2(w_2 * g_1(w_1 * x)), y^*)$$

3. Now we will work through parts of the backward pass, incrementally. Use the chain rule to derive $\frac{\partial Loss}{\partial w_2}$. Write your expression as a product of partial derivatives at each node: i.e. the partial derivative of the node's output with respect to its inputs. (Hint: the series of expressions you wrote in part 1 will be helpful; you may use any of those variables.)

$$\frac{\partial Loss}{\partial w_2} = \frac{\partial Loss}{\partial a_2} \frac{\partial a_2}{\partial z_2} \frac{\partial z_2}{\partial w_2}$$

4. Suppose the loss function is quadratic, $Loss(a_2, y^*) = \frac{1}{2}(a_2 - y^*)^2$, and g_1 and g_2 are both sigmoid functions $g(z) = \frac{1}{1+e^{-z}}$ (note: it's typically better to use a different type of loss, *cross-entropy*, for classification problems, but we'll use this to make the math easier).

Using the chain rule from Part 3, and the fact that $\frac{\partial g(z)}{\partial z} = g(z)(1 - g(z))$ for the sigmoid function, write $\frac{\partial Loss}{\partial w_2}$ in terms of the values from the forward pass, y^* , a_1 , and a_2 :

First we'll compute the partial derivatives at each node:

$$\begin{aligned}\frac{\partial Loss}{\partial a_2} &= (a_2 - y^*) \\ \frac{\partial a_2}{\partial z_2} &= \frac{\partial g_2(z_2)}{\partial z_2} = g_2(z_2)(1 - g_2(z_2)) = a_2(1 - a_2) \\ \frac{\partial z_2}{\partial w_2} &= a_1\end{aligned}$$

Now we can plug into the chain rule from part 3:

$$\begin{aligned}\frac{\partial Loss}{\partial w_2} &= \frac{\partial Loss}{\partial a_2} \frac{\partial a_2}{\partial z_2} \frac{\partial z_2}{\partial w_2} \\ &= (a_2 - y^*) * a_2(1 - a_2) * a_1\end{aligned}$$

5. Now use the chain rule to derive $\frac{\partial Loss}{\partial w_1}$ as a product of partial derivatives at each node used in the chain rule:

$$\frac{\partial Loss}{\partial w_1} = \frac{\partial Loss}{\partial a_2} \frac{\partial a_2}{\partial z_2} \frac{\partial z_2}{\partial a_1} \frac{\partial a_1}{\partial z_1} \frac{\partial z_1}{\partial w_1}$$

6. Finally, write $\frac{\partial Loss}{\partial w_1}$ in terms of x, y^*, w_i, a_i, z_i : The partial derivatives at each node (in addition to the ones we computed in Part 4) are:

$$\begin{aligned}\frac{\partial z_2}{\partial a_1} &= w_2 \\ \frac{\partial a_1}{\partial z_1} &= \frac{\partial g_1(z_1)}{\partial z_1} = g_1(z_1)(1 - g_1(z_1)) = a_1(1 - a_1) \\ \frac{\partial z_1}{\partial a_1} &= x\end{aligned}$$

Plugging into the chain rule from Part 5 gives:

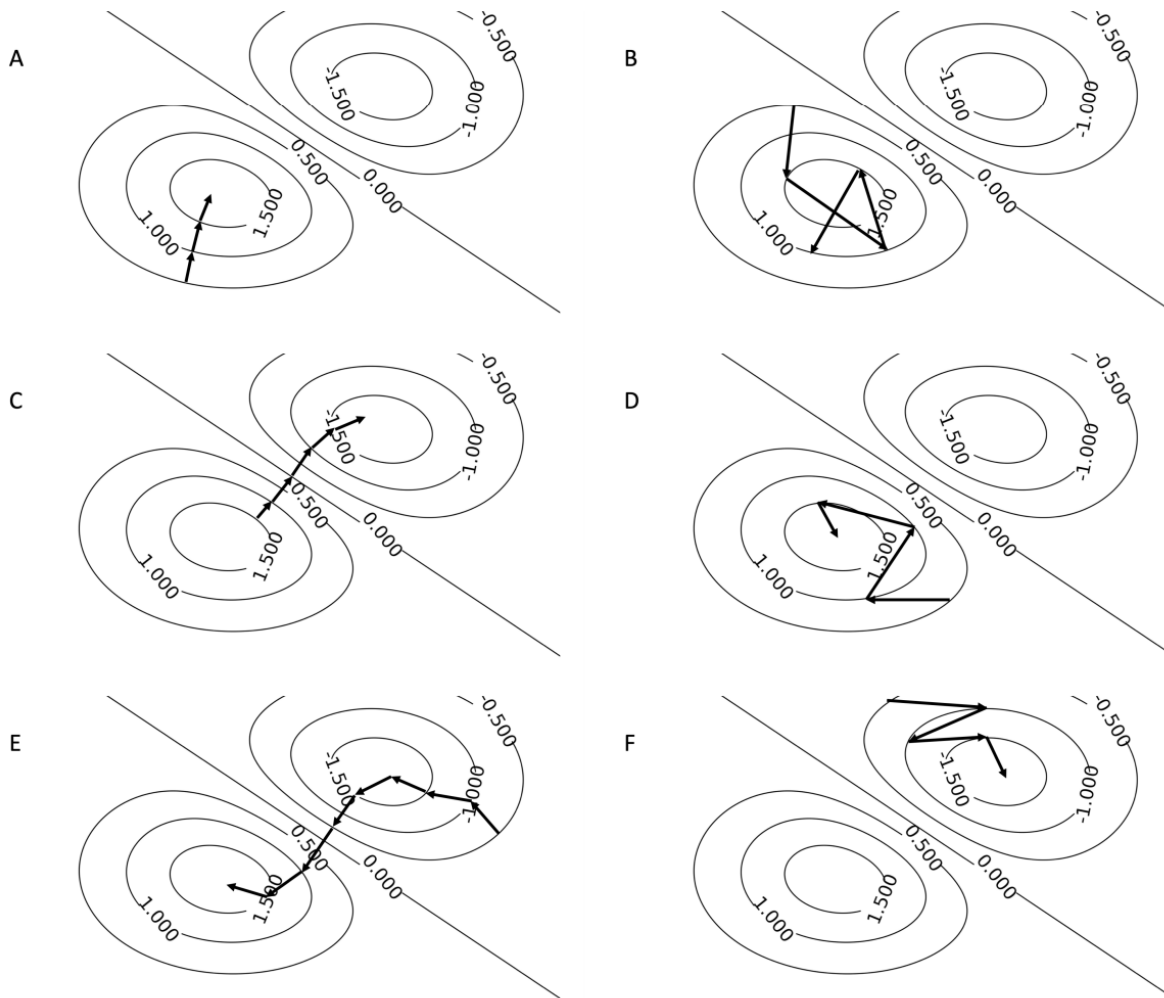
$$\begin{aligned}\frac{\partial Loss}{\partial w_1} &= \frac{\partial Loss}{\partial a_2} \frac{\partial a_2}{\partial z_2} \frac{\partial z_2}{\partial a_1} \frac{\partial a_1}{\partial z_1} \frac{\partial z_1}{\partial w_1} \\ &= (a_2 - y^*) * a_2(1 - a_2) * w_2 * a_1(1 - a_1) * x\end{aligned}$$

7. What is the gradient descent update for w_1 with step-size α in terms of the values computed above?

$$w_1 \leftarrow w_1 - \alpha(a_2 - y^*) * a_2(1 - a_2) * w_2 * a_1(1 - a_1) * x$$

Q3. Gradient Ascent Trajectory

(a) Which of the following paths is a feasible trajectory for the gradient ascent algorithm?

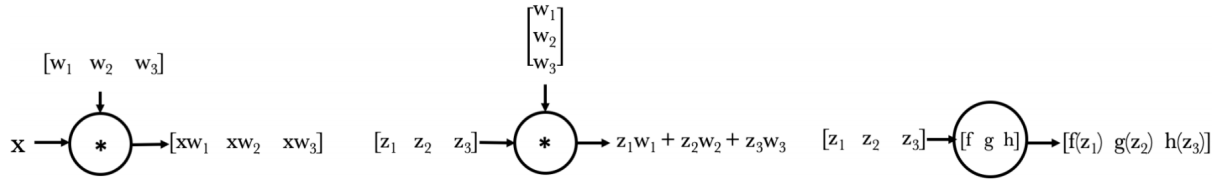


- A
 B
 C
 D
 E
 F

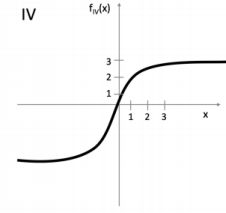
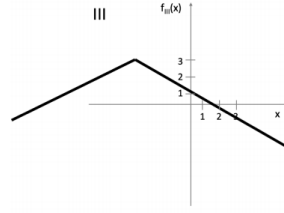
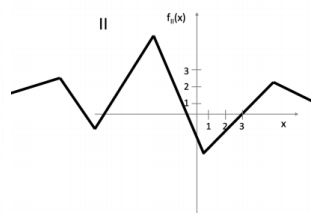
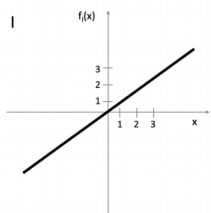
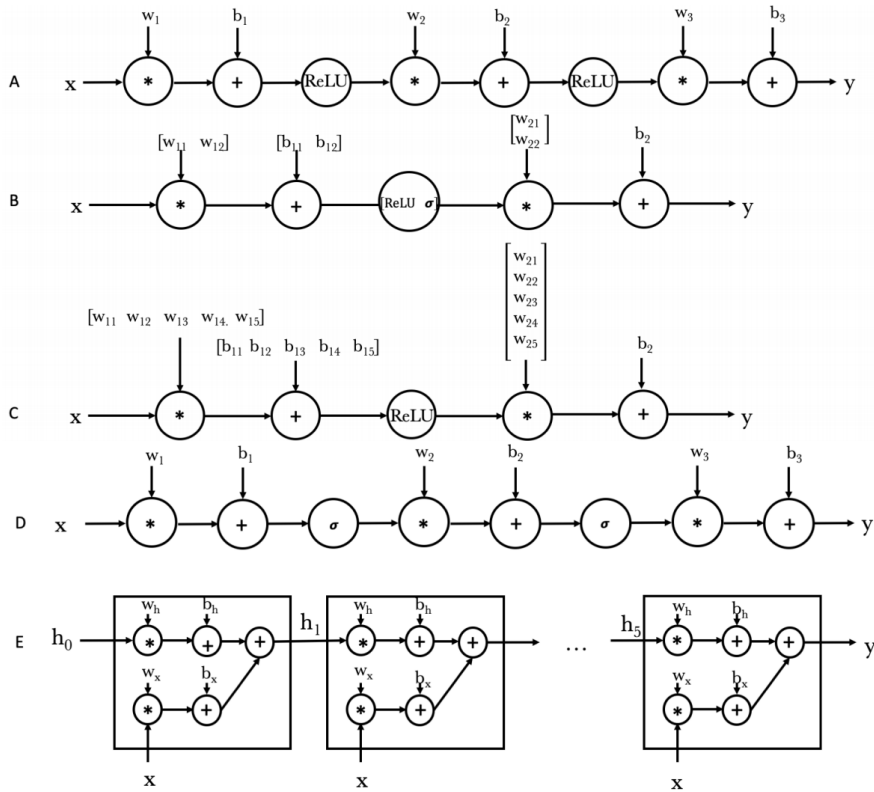
A is a gradient ascent path since the gradient lines are orthogonal to the contours and the point towards the maximum. B is also a gradient ascent path with a high learning rate. C is not because the path is going towards the minimum instead of the maximum. D is not a gradient ascent path since the gradient is not orthogonal to the contour lines. E is not a gradient ascent path since it starts going towards the minimum. F is not since it goes towards the minimum and the gradients are not orthogonal to the contour lines.

Q4. Neural Networks Representation

(a) We are given the following 5 neural networks (NN) architectures. The operation $*$ represents the matrix multiplication operation, $[w_{i1} \dots w_{ik}]$ and $[b_{i1} \dots b_{ik}]$ represents the weights and the biases of the NN, the orientation (vertical and horizontal) is just for consistency in the operations. The term $[\text{ReLU } \sigma]$ in B means applying a ReLU activation to the first element of the vector and a sigmoid (σ) activation to the second element. These operations are depicted in the following figures:



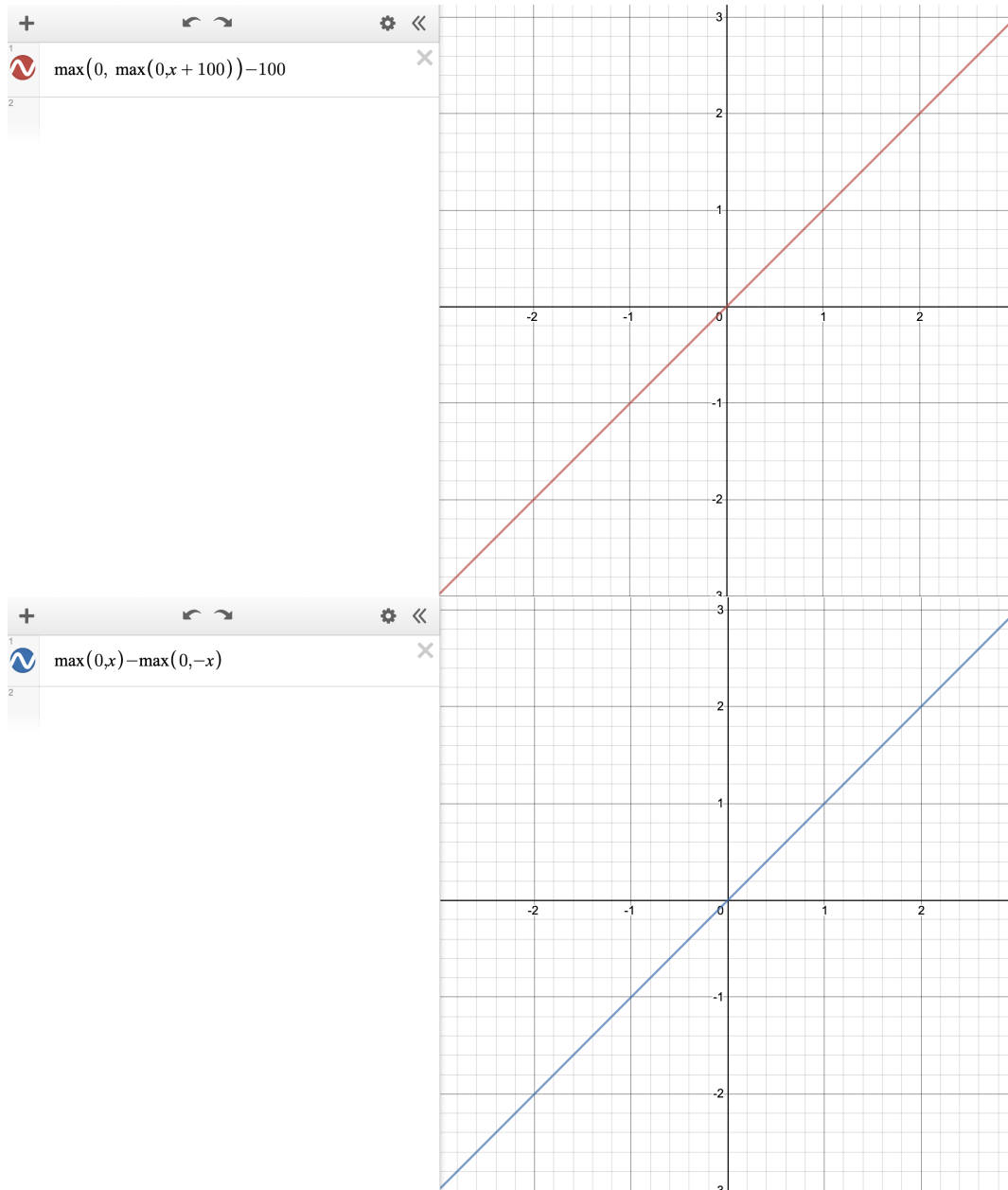
Which of the following neural networks can represent each function?



(i) $f_{\text{I}}(x)$:

A B C D E

A can represent it by using sufficiently high b_1 to shift the input domain to where ReLU appears the same as the identity function then a sufficiently negative b_3 to shift the function back to the identity, i.e. with $\max(0, \max(0, x + 100)) - 100$. *B* can do so similarly by setting the parameter for σ to be 0 so we are essentially back to the ReLU case. *C* can represent it by having $w_{11} = 1, w_{12} = -1, w_{21} = 1, w_{22} = -1$ and the rest of the parameters being 0, i.e. $\max(0, x) - \max(0, -x)$. *D* cannot because of sigmoid's smooth and non-linear activations. *E* is a linear function and therefore can represent the identity.



(ii) $f_{\text{II}}(x)$:

A B C D E

This is a piecewise linear function with 6 pieces. As a result you can represent it by linearly combining 5 ReLU functions. The only possible graph that can represent this function is C .

(iii) $f_{\text{III}}(x)$:

A B C D E

This is a piecewise linear function with 2 pieces. This can be obtained by linearly combining two ReLU functions. The only possible solution is then C . Note that A cannot represent this function since the last ReLU of the network would result in a flat semi-line.

(iv) $f_{\text{IV}}(x)$:

A B C D E

A, C, E cannot represent any non-linear or nonpiecewise linear functions.