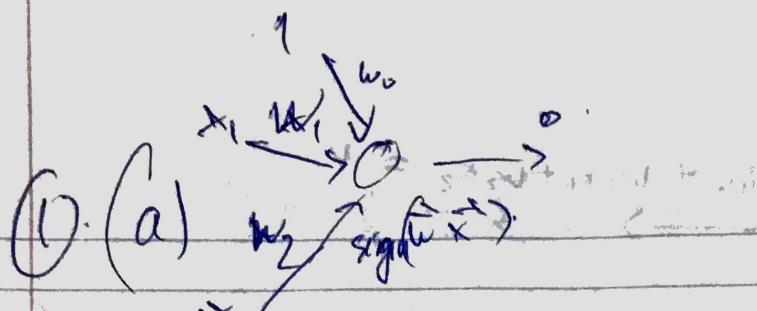


83/21



$$\Phi = y = \text{sign} \{ \vec{w} \cdot \vec{x} \} = \text{sign} \{ \sum_{i=1}^2 w_i x_i \} = 0 \text{ (not zero.)}$$

$$(b) b = x_0 \cdot w_0 = 1 \cdot (0.1) = 0.1 \quad (\text{Bias})$$

$$\begin{array}{c} x_0 = 1 \\ x_1 \\ x_2 \end{array} \xrightarrow{\begin{array}{l} w_0 \\ w_1 \\ w_2 \end{array}} 0 \rightarrow$$

$$y = \text{sign} \{ w_0 + w_1 x_1 + w_2 x_2 \} = 0$$

$$= \text{sign} \{ 0.1 + 0.2(2) + (-0.3)(-3) \}$$

$$(c) \hat{y} = \text{sign} \{ 0.1 + 0.2(2) + (-0.3)(-3) \}$$

$$= \text{sign} \{ -8.5 \} = -1.$$

$$(d) L = (y - \hat{y})^2 = (1 - (-1))^2 = 2^2 = 4.$$

$$(e) \vec{w} \leftarrow \vec{w} + \alpha (y - \hat{y}) \vec{x}$$

$$w = (0.1, 0.2, -0.3) + \alpha (4) (1, 2, 3)$$

$$= (0.1, 0.2, -0.3) + (4\alpha, 8\alpha, 12\alpha)$$

$$= (0.1 + 4\alpha, 0.2 + 8\alpha, -0.3 + 12\alpha)$$

(2) (a)
$$x_1 \xrightarrow{w_1} \textcircled{B} \xrightarrow{w_0 + w_1 x_1 + w_2 x_2 = v}$$

$$x_2 \xrightarrow{w_2}$$

$$v = w_0 x_0 + w_0(1) = w_0 \text{ (bias)}$$

(b) $v = v$ (Identity)

$$\Phi(w_0 + w_1 x_1 + w_2 x_2) = w_0 + w_1 x_1 + w_2 x_2.$$

(b) $\frac{\partial L}{\partial x_1} = \frac{\partial L}{\partial v} \cdot \frac{\partial v}{\partial x_1} = \cancel{(5)}(5)(0.2) = 1.$

$$\frac{\partial v}{\partial x_1} = w_1.$$

(c) $\frac{\partial L}{\partial x_2} = \frac{\partial L}{\partial v} \cdot \frac{\partial v}{\partial x_2} = \cancel{(5)}(-0.3) = -1.5$

$$\frac{\partial v}{\partial x_2} = w_2.$$

$$\frac{\partial L}{\partial x_0} = \frac{\partial L}{\partial v} \cdot \frac{\partial v}{\partial x_0} = \cancel{(5)}(0) = 0$$

$$\frac{\partial v}{\partial x_0} = 0 \quad (1) \quad (0) = 0 \quad (0) = 0$$

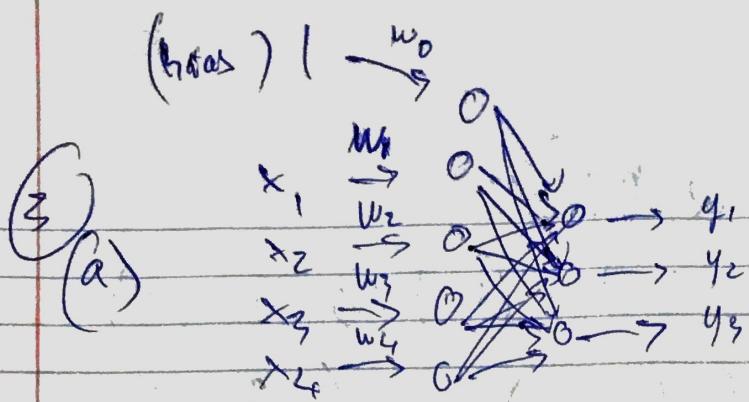
$$\Rightarrow \nabla_{w^T} L = \begin{pmatrix} \frac{\partial L}{\partial w_0} \\ \frac{\partial L}{\partial w_1} \\ \frac{\partial L}{\partial w_2} \end{pmatrix} = \begin{pmatrix} 0 \\ 1 \\ -1.5 \end{pmatrix}$$

(2)

(d) $\vec{w} \leftarrow \vec{w} - \alpha \nabla_{\vec{w}} L$

$$w = (0.1, 0.2, -0.3) - (0.1)(0, 1, -1.5)$$

$$= \begin{pmatrix} 0.1 - 0.1(0) \\ 0.2 - 0.1(1) \\ -0.3 - 0.1(-1.5) \end{pmatrix} = \begin{pmatrix} 0 \\ 0.1 \\ -2.45 \end{pmatrix}$$



4 inputs because there are 4 features

(b) 3 outputs because there are 3 different classes

(c) 8 weights (including bias)

(d) Multinomial logistic Regression is a Softmax classifier.
 \Rightarrow use softmax activation and
 use cross-entropy as loss function.

(e) Given \vec{x} , ~~we~~^{The} multinomial logistic regression
 will give us the outputs that contain likelihood or
 probabilities of each class - Based on the probability,
 likelihood/probability
 the prediction is the highest or most likely class
 to occur.

(4) (a) if α is too small, the learning rate becomes slow. \rightarrow Takes longer time to train if we want to get high accuracy.

(b) # of neuron in hidden layer is low \rightarrow the accuracy is low because model underfits training set.

(c) # of epochs is low \rightarrow The ~~total~~ learning may not completed or converge \rightarrow the accuracy is low and, high loss.

(d) if λ is ~~less~~^{small}, then model cannot penalize the weights. Some of weights ~~may~~ ^{should not have high impact/influence} on some features. The regulariser helps us to lower this influence. ~~so it has~~

(e) mini-batch size is small \rightarrow The training cannot learn fast. The group of update can be trained at the same \Rightarrow learn slow and may effect accuracy.

(5) (iii) Some of potential causes

- Use too much dense layers

(a) This can have negative effect like overfitting the training set.

(b) They can use fewer ~~sets~~ of dense layers.

(c) if they change as suggested in (b), the accuracy and loss may get improved.

- Use ~~inappropriate~~ activation function and loss function.

(a) This will cause the wrong prediction or unexpected results.

(b) They can change ~~sets~~ to different activation functions and loss functions

(c) if they ~~choose~~ apply the change, the outputs can be better.

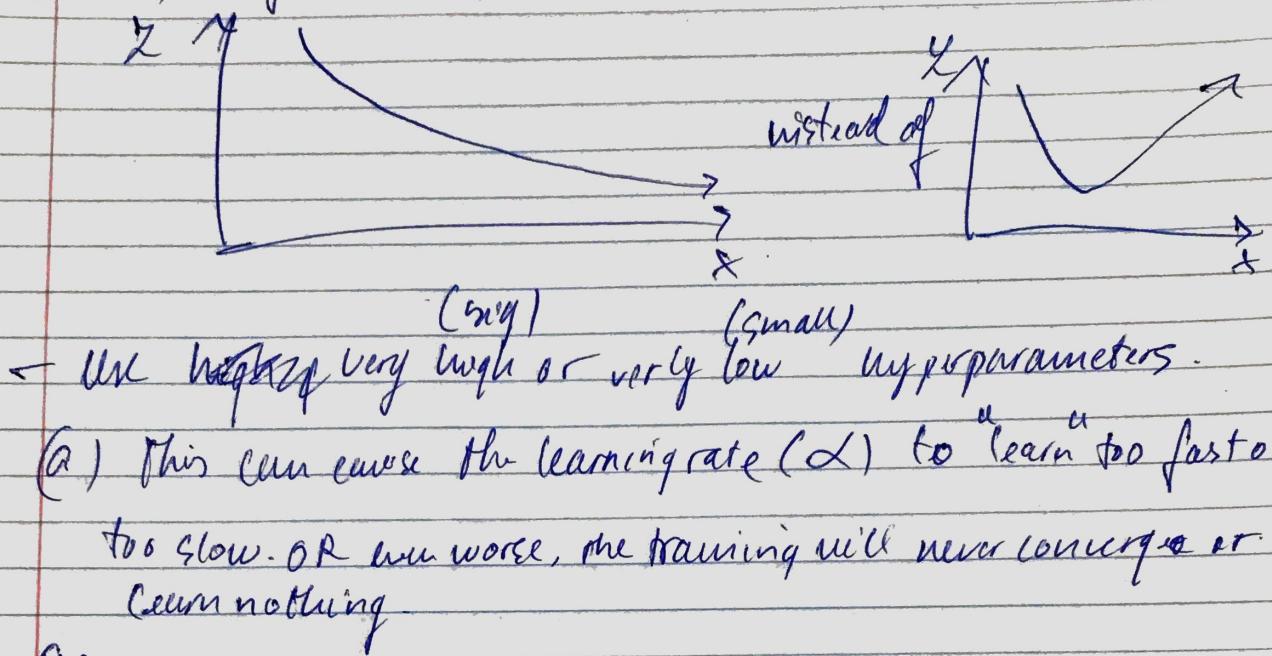
- Training data set may be too small

(a) This can ~~lead to~~ ~~not~~ make the model (not get information about, learn the pattern and leads to approximation or overfitting).

(b) They can change to larger training set.

(c) if they apply the change, they would expect a better result. (e.g. higher accuracy, lower loss)

⑤ (c) They would expect something like this on graph.



← Use ~~high~~ very high or very low hyperparameters.

(a) This can cause the learning rate (α) to "learn" too fast or too slow. OR even worse, the training will never converge or learn nothing.

(b) they can tune hyperparams to the appropriate values

(c) They would expect the see the loss as ~~faster~~ picture ^{above} on the left.