

Module 4.3: Output Functions and Loss Functions

We need to answer two questions

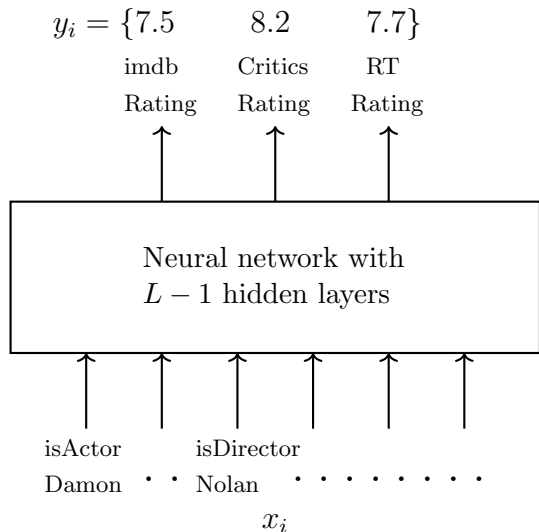
- How to choose the loss function $\mathcal{L}(\theta)$?
- How to compute $\nabla\theta$ which is composed of
 $\nabla W_1, \nabla W_2, \dots, \nabla W_{L-1} \in \mathbb{R}^{n \times n}, \nabla W_L \in \mathbb{R}^{n \times k},$
 $\nabla b_1, \nabla b_2, \dots, \nabla b_{L-1} \in \mathbb{R}^n$ and $\nabla b_L \in \mathbb{R}^k$?

We need to answer two questions

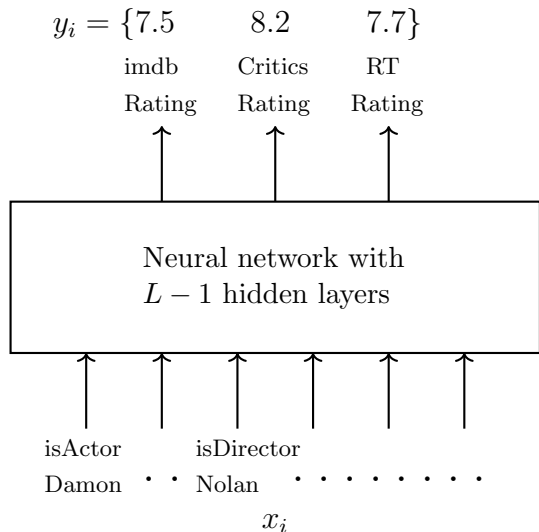
- How to choose the loss function $\mathcal{L}(\theta)$?
- How to compute $\nabla\theta$ which is composed of $\nabla W_1, \nabla W_2, \dots, \nabla W_{L-1} \in \mathbb{R}^{n \times n}, \nabla W_L \in \mathbb{R}^{n \times k}, \nabla b_1, \nabla b_2, \dots, \nabla b_{L-1} \in \mathbb{R}^n$ and $\nabla b_L \in \mathbb{R}^k$?

- The choice of loss function depends on the problem at hand

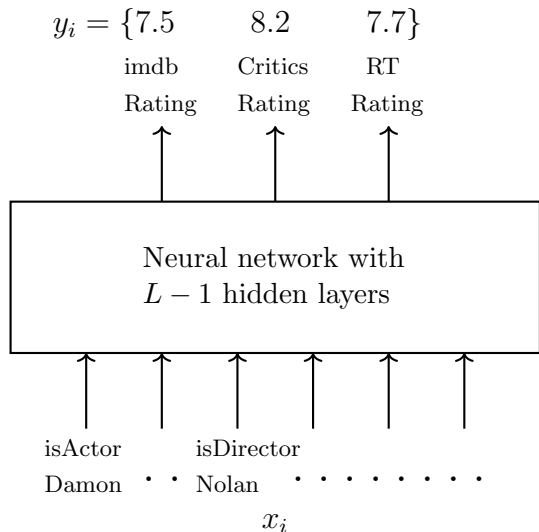
- The choice of loss function depends on the problem at hand
- We will illustrate this with the help of two examples



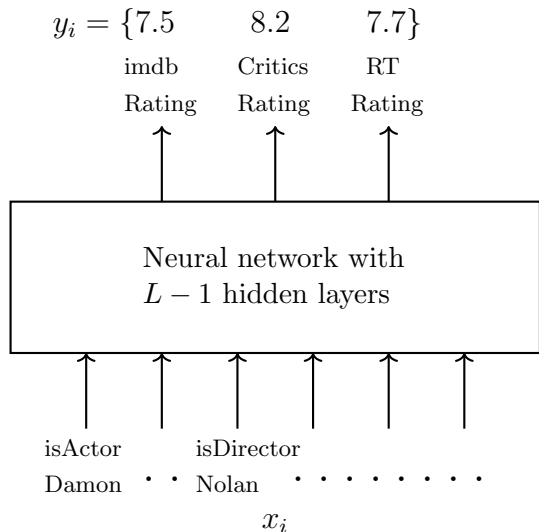
- The choice of loss function depends on the problem at hand
- We will illustrate this with the help of two examples
- Consider our movie example again but this time we are interested in predicting ratings



- The choice of loss function depends on the problem at hand
- We will illustrate this with the help of two examples
- Consider our movie example again but this time we are interested in predicting ratings
- Here $y_i \in \mathbb{R}^3$



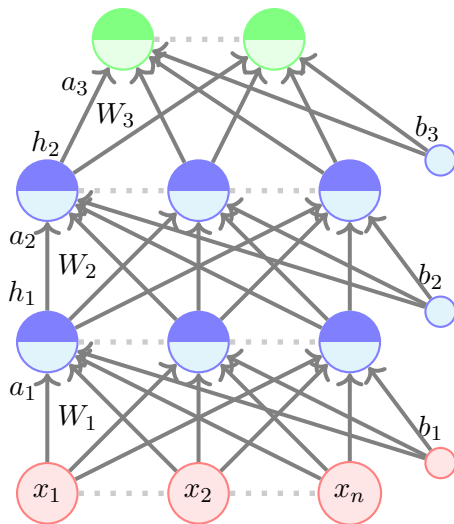
- The choice of loss function depends on the problem at hand
- We will illustrate this with the help of two examples
- Consider our movie example again but this time we are interested in predicting ratings
- Here $y_i \in \mathbb{R}^3$
- The loss function should capture how much \hat{y}_i deviates from y_i



- The choice of loss function depends on the problem at hand
- We will illustrate this with the help of two examples
- Consider our movie example again but this time we are interested in predicting ratings
- Here $y_i \in \mathbb{R}^3$
- The loss function should capture how much \hat{y}_i deviates from y_i
- If $y_i \in \mathbb{R}^n$ then the squared error loss can capture this deviation

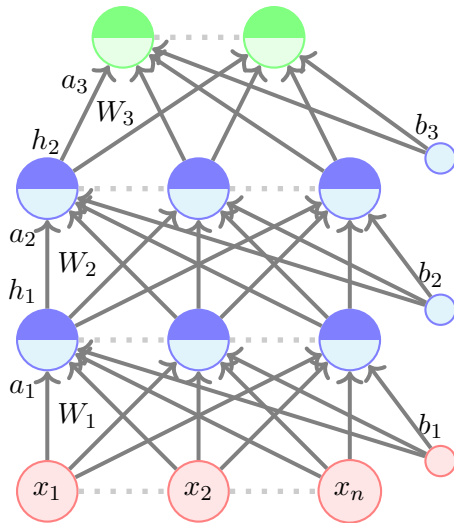
$$\mathcal{L}(\theta) = \frac{1}{N} \sum_{i=1}^N \sum_{j=1}^3 (\hat{y}_{ij} - y_{ij})^2$$

$$h_L = \hat{y} = f(x)$$



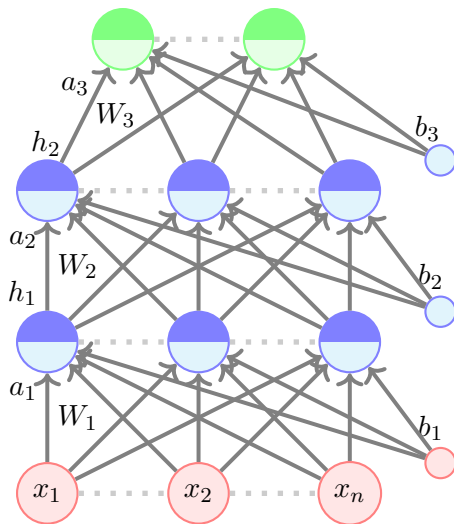
- A related question: What should the output function ‘ O ’ be if $y_i \in \mathbb{R}$?

$$h_L = \hat{y} = f(x)$$



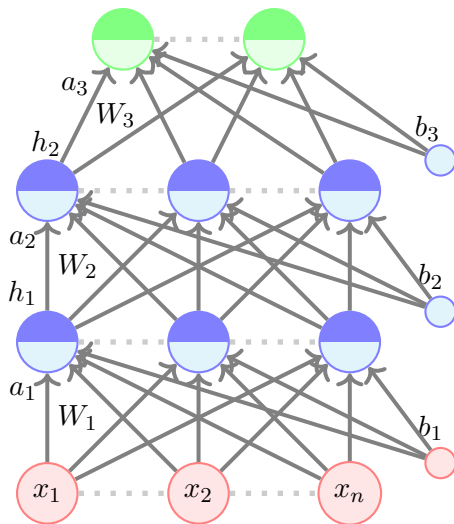
- A related question: What should the output function ‘ O ’ be if $y_i \in \mathbb{R}$?
- More specifically, can it be the logistic function?

$$h_L = \hat{y} = f(x)$$



- A related question: What should the output function ‘ O ’ be if $y_i \in \mathbb{R}$?
- More specifically, can it be the logistic function?
- No, because it restricts \hat{y}_i to a value between 0 & 1 but we want $\hat{y}_i \in \mathbb{R}$

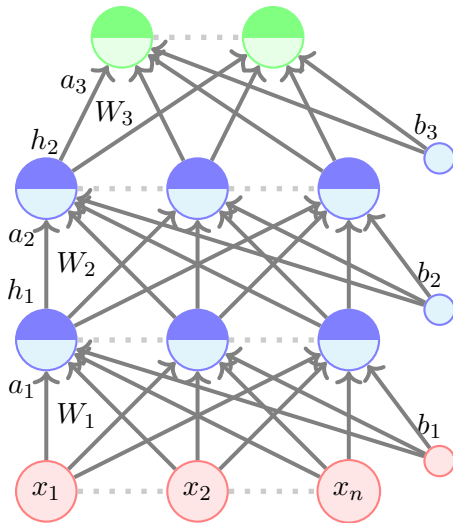
$$h_L = \hat{y} = f(x)$$



- A related question: What should the output function ‘ O ’ be if $y_i \in \mathbb{R}$?
- More specifically, can it be the logistic function?
- No, because it restricts \hat{y}_i to a value between 0 & 1 but we want $\hat{y}_i \in \mathbb{R}$
- So, in such cases it makes sense to have ‘ O ’ as identity function

$$\begin{aligned} f(x) &= h_L = O(a_L) \\ &= a_L = W_L h_{L-1} + b_L \end{aligned}$$

$$h_L = \hat{y} = f(x)$$



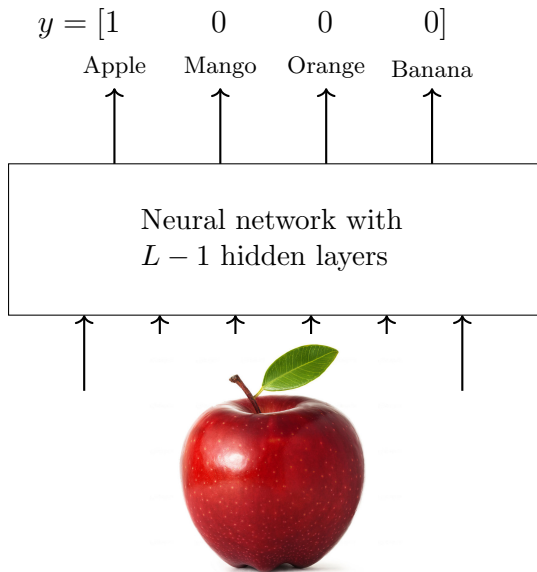
- A related question: What should the output function ‘ O ’ be if $y_i \in \mathbb{R}$?
- More specifically, can it be the logistic function?
- No, because it restricts \hat{y}_i to a value between 0 & 1 but we want $\hat{y}_i \in \mathbb{R}$
- So, in such cases it makes sense to have ‘ O ’ as identity function

$$\begin{aligned} f(x) &= h_L = O(a_L) \\ &= a_L = W_L h_{L-1} + b_L \end{aligned}$$

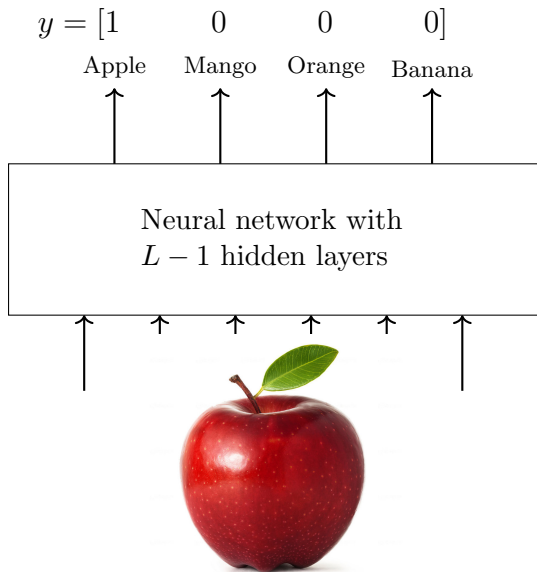
- $\hat{y}_i = f(x_i)$ is no longer bounded between 0 and 1

Intentionally left blank

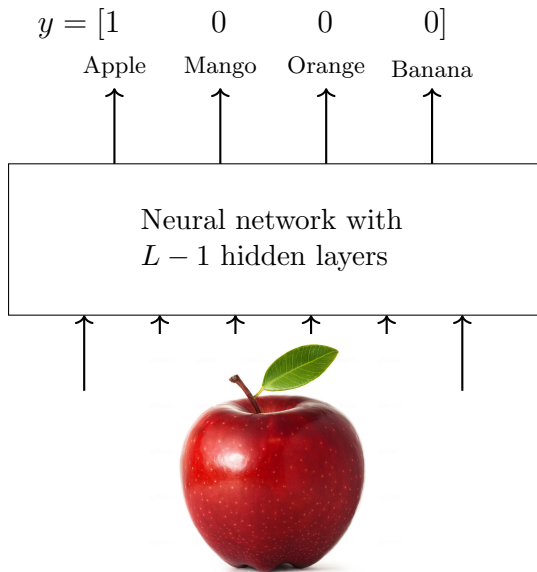
Intentionally left blank



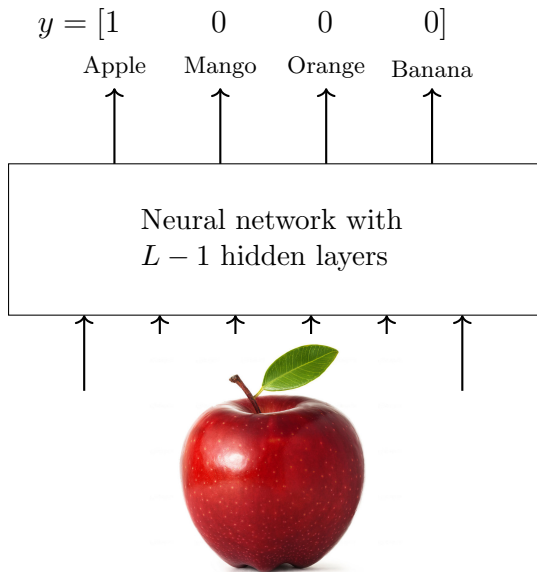
- Now let us consider another problem for which a different loss function would be appropriate



- Now let us consider another problem for which a different loss function would be appropriate
- Suppose we want to classify an image into 1 of k classes

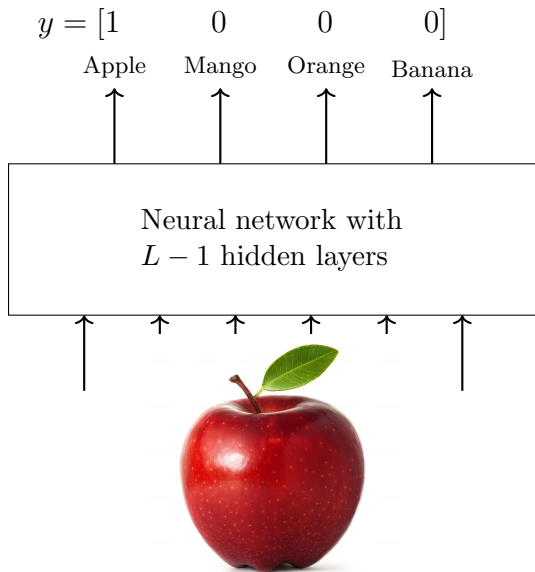


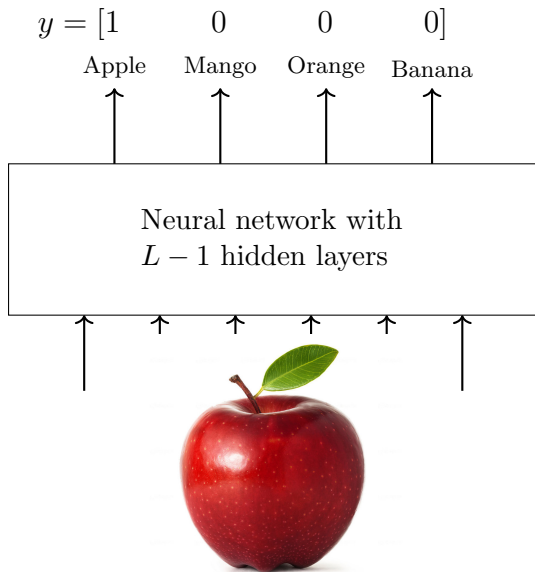
- Now let us consider another problem for which a different loss function would be appropriate
- Suppose we want to classify an image into 1 of k classes
- Here again we could use the squared error loss to capture the deviation



- Now let us consider another problem for which a different loss function would be appropriate
- Suppose we want to classify an image into 1 of k classes
- Here again we could use the squared error loss to capture the deviation
- But can you think of a better function?

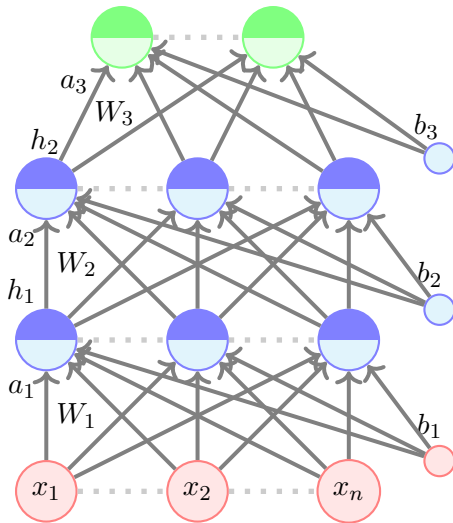
- Notice that y is a probability distribution





- Notice that y is a probability distribution
- Therefore we should also ensure that \hat{y} is a probability distribution

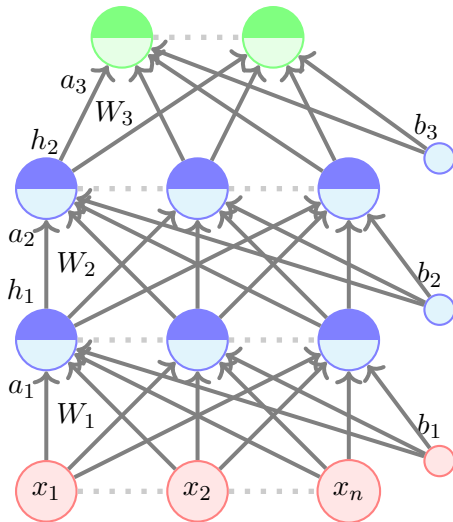
$$h_L = \hat{y} = f(x)$$



- Notice that y is a probability distribution
- Therefore we should also ensure that \hat{y} is a probability distribution
- What choice of the output activation ‘ O ’ will ensure this ?

$$a_L = W_L h_{L-1} + b_L$$

$$h_L = \hat{y} = f(x)$$



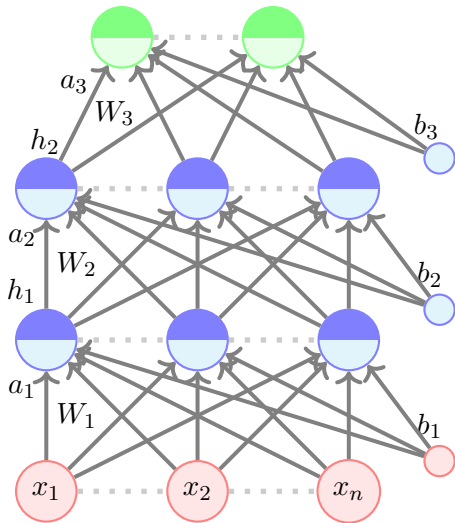
- Notice that y is a probability distribution
- Therefore we should also ensure that \hat{y} is a probability distribution
- What choice of the output activation ‘ O ’ will ensure this ?

$$a_L = W_L h_{L-1} + b_L$$

$$f(x)_j = O(a_L)_j = \frac{e^{a_{L,j}}}{\sum_{j'=1}^k e^{a_{L,j'}}}$$

$O(a_L)_j$ is the j^{th} element of \hat{y}
 $a_{L,j}$ is the j^{th} element of the vector a_L .

$$h_L = \hat{y} = f(x)$$



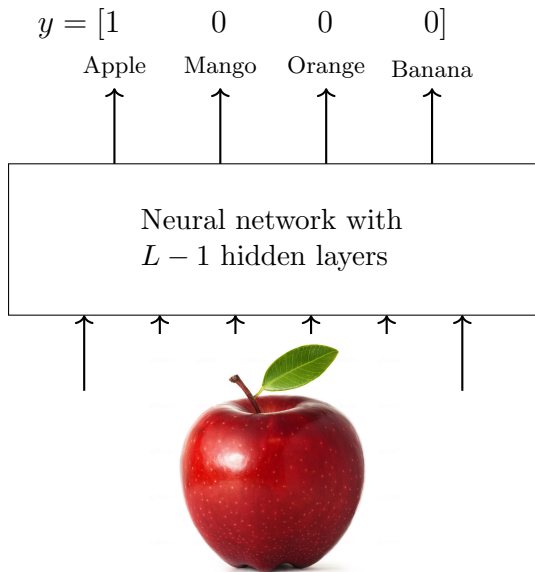
- Notice that y is a probability distribution
- Therefore we should also ensure that \hat{y} is a probability distribution
- What choice of the output activation ‘ O ’ will ensure this ?

$$a_L = W_L h_{L-1} + b_L$$

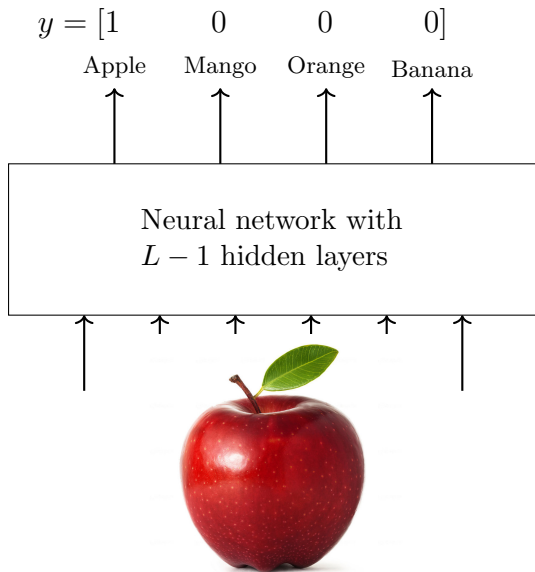
$$f(x)_j = O(a_L)_j = \frac{e^{a_{L,j}}}{\sum_{j'=1}^k e^{a_{L,j'}}}$$

$O(a_L)_j$ is the j^{th} element of \hat{y}
 $a_{L,j}$ is the j^{th} element of the vector a_L .

- This function is called the *softmax* function

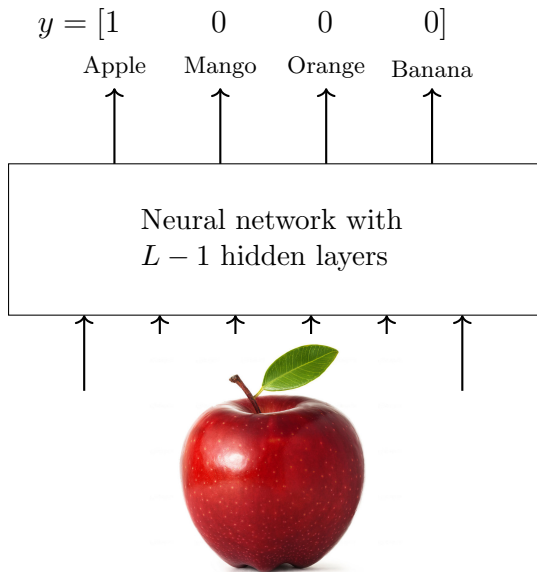


- Now that we have ensured that both y & \hat{y} are probability distributions can you think of a function which captures the difference between them?



- Now that we have ensured that both y & \hat{y} are probability distributions can you think of a function which captures the difference between them?
- Cross-entropy

$$\mathcal{L}(\theta) = - \sum_{c=1}^k y_c \log \hat{y}_c$$



- Now that we have ensured that both y & \hat{y} are probability distributions can you think of a function which captures the difference between them?
- Cross-entropy

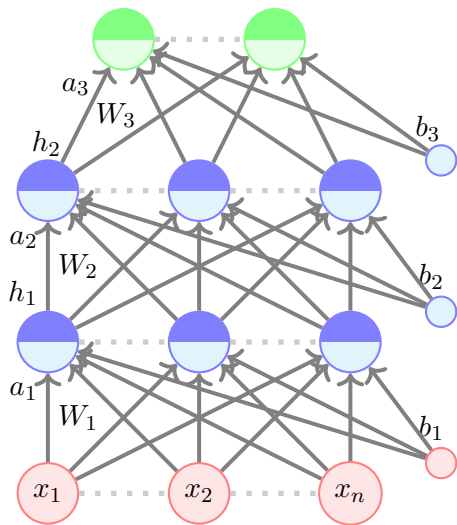
$$\mathcal{L}(\theta) = - \sum_{c=1}^k y_c \log \hat{y}_c$$

- Notice that

$$y_c = \begin{cases} 1 & \text{if } c = \ell \text{ (the true class label)} \\ 0 & \text{otherwise} \end{cases}$$

$$\mathcal{L}(\theta) = -\log \hat{y}_\ell$$

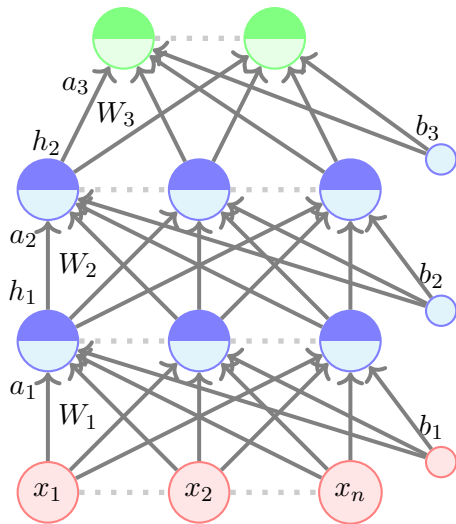
$$h_L = \hat{y} = f(x)$$



- So, for classification problem (where you have to choose 1 of K classes), we use the following objective function

$$\begin{aligned} & \underset{\theta}{\text{minimize}} && \mathcal{L}_{\theta} = -\log \hat{y}_{\ell} \\ \text{or} & && \underset{\theta}{\text{maximize}} && -\mathcal{L}_{\theta} = \log \hat{y}_{\ell} \end{aligned}$$

$$h_L = \hat{y} = f(x)$$

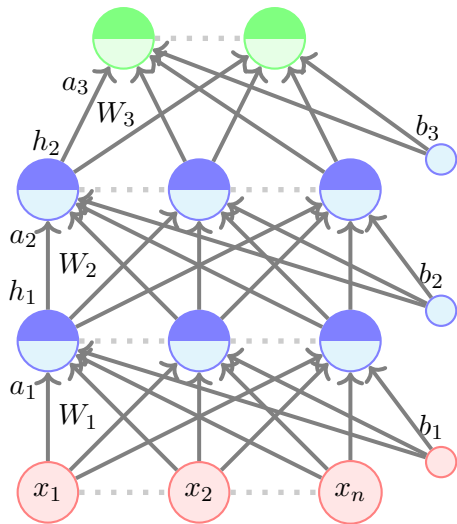


- So, for classification problem (where you have to choose 1 of K classes), we use the following objective function

$$\begin{aligned} & \underset{\theta}{\text{minimize}} && \mathcal{L}_{\theta} = -\log \hat{y}_{\ell} \\ \text{or} & && \underset{\theta}{\text{maximize}} && -\mathcal{L}_{\theta} = \log \hat{y}_{\ell} \end{aligned}$$

- But wait!
Is \hat{y}_{ℓ} a function of $\theta = [W_1, W_2, \dots, W_L, b_1, b_2, \dots, b_L]$?

$$h_L = \hat{y} = f(x)$$

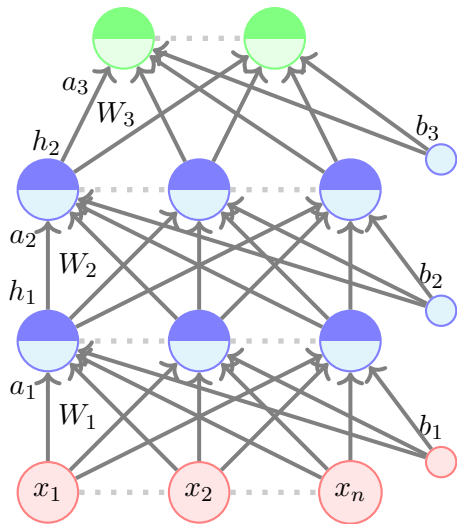


- So, for classification problem (where you have to choose 1 of K classes), we use the following objective function

$$\begin{aligned} & \underset{\theta}{\text{minimize}} && \mathcal{L}_{\theta} = -\log \hat{y}_{\ell} \\ \text{or} &&& \underset{\theta}{\text{maximize}} && -\mathcal{L}_{\theta} = \log \hat{y}_{\ell} \end{aligned}$$

- But wait!
Is \hat{y}_{ℓ} a function of $\theta = [W_1, W_2, \dots, W_L, b_1, b_2, \dots, b_L]$?
- Yes, it is indeed a function of θ
$$\hat{y}_{\ell} = [O(W^3 g(W^2 g(W^1 x + b_1) + b_2) + b_3)]_{\ell}$$

$$h_L = \hat{y} = f(x)$$



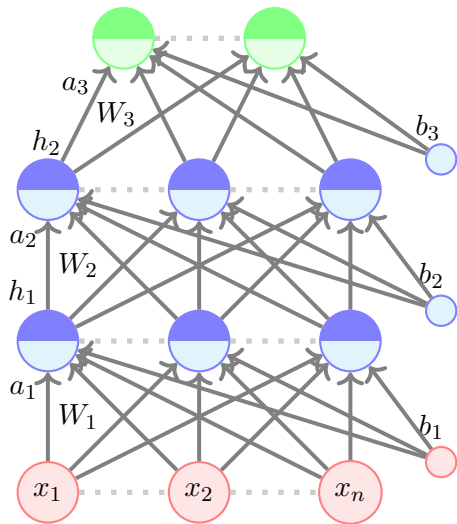
- So, for classification problem (where you have to choose 1 of K classes), we use the following objective function

$$\begin{aligned} & \underset{\theta}{\text{minimize}} && \mathcal{L}_{\theta} = -\log \hat{y}_{\ell} \\ \text{or} &&& \underset{\theta}{\text{maximize}} && -\mathcal{L}_{\theta} = \log \hat{y}_{\ell} \end{aligned}$$

- But wait!
Is \hat{y}_{ℓ} a function of $\theta = [W_1, W_2, \dots, W_L, b_1, b_2, \dots, b_L]$?
- Yes, it is indeed a function of θ

$$\hat{y}_{\ell} = [O(W^3 g(W^2 g(W^1 x + b_1) + b_2) + b_3)]_{\ell}$$
- What does \hat{y}_{ℓ} encode?

$$h_L = \hat{y} = f(x)$$



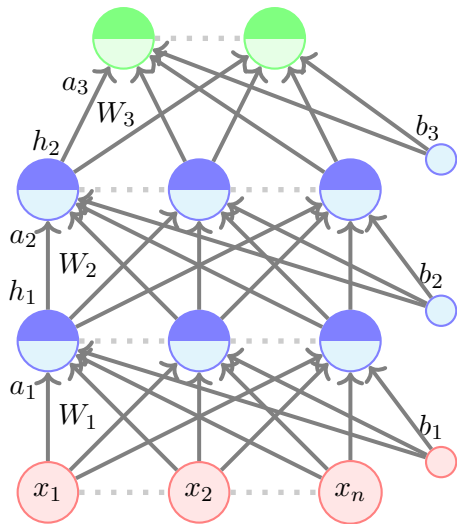
- So, for classification problem (where you have to choose 1 of K classes), we use the following objective function

$$\begin{aligned} & \underset{\theta}{\text{minimize}} && \mathcal{L}_\theta = -\log \hat{y}_\ell \\ \text{or} &&& \underset{\theta}{\text{maximize}} && -\mathcal{L}_\theta = \log \hat{y}_\ell \end{aligned}$$

- But wait!
Is \hat{y}_ℓ a function of $\theta = [W_1, W_2, \dots, W_L, b_1, b_2, \dots, b_L]$?
- Yes, it is indeed a function of θ

$$\hat{y}_\ell = [O(W^3 g(W^2 g(W^1 x + b_1) + b_2) + b_3)]_\ell$$
- What does \hat{y}_ℓ encode?
- It is the probability that x belongs to the ℓ^{th} class (bring it as close to 1).

$$h_L = \hat{y} = f(x)$$



- So, for classification problem (where you have to choose 1 of K classes), we use the following objective function

$$\begin{aligned} & \underset{\theta}{\text{minimize}} && \mathcal{L}_\theta = -\log \hat{y}_\ell \\ \text{or} &&& \underset{\theta}{\text{maximize}} && -\mathcal{L}_\theta = \log \hat{y}_\ell \end{aligned}$$

- But wait!
Is \hat{y}_ℓ a function of $\theta = [W_1, W_2, \dots, W_L, b_1, b_2, \dots, b_L]$?
- Yes, it is indeed a function of θ

$$\hat{y}_\ell = [O(W^3 g(W^2 g(W^1 x + b_1) + b_2) + b_3)]_\ell$$
- What does \hat{y}_ℓ encode?
- It is the probability that x belongs to the ℓ^{th} class (bring it as close to 1).
- $\log \hat{y}_\ell$ is called the *log-likelihood* of the data.

| | Outputs | |
|-------------------|----------------|---------------|
| | Real Values | Probabilities |
| Output Activation | | |
| Loss Function | | |

| | | |
|-------------------|----------------|---------------|
| | Outputs | |
| | Real Values | Probabilities |
| Output Activation | Linear | |
| Loss Function | | |

| | Outputs | |
|-------------------|-------------|---------------|
| | Real Values | Probabilities |
| Output Activation | Linear | Softmax |
| Loss Function | | |

| | Outputs | |
|-------------------|----------------|---------------|
| | Real Values | Probabilities |
| Output Activation | Linear | Softmax |
| Loss Function | Squared Error | |

| | Outputs | |
|-------------------|----------------|---------------|
| | Real Values | Probabilities |
| Output Activation | Linear | Softmax |
| Loss Function | Squared Error | Cross Entropy |

| | Outputs | |
|-------------------|---------------|---------------|
| | Real Values | Probabilities |
| Output Activation | Linear | Softmax |
| Loss Function | Squared Error | Cross Entropy |

- Of course, there could be other loss functions depending on the problem at hand but the two loss functions that we just saw are encountered very often

| | Outputs | |
|-------------------|---------------|---------------|
| | Real Values | Probabilities |
| Output Activation | Linear | Softmax |
| Loss Function | Squared Error | Cross Entropy |

- Of course, there could be other loss functions depending on the problem at hand but the two loss functions that we just saw are encountered very often
- For the rest of this lecture we will focus on the case where the output activation is a softmax function and the loss function is cross entropy