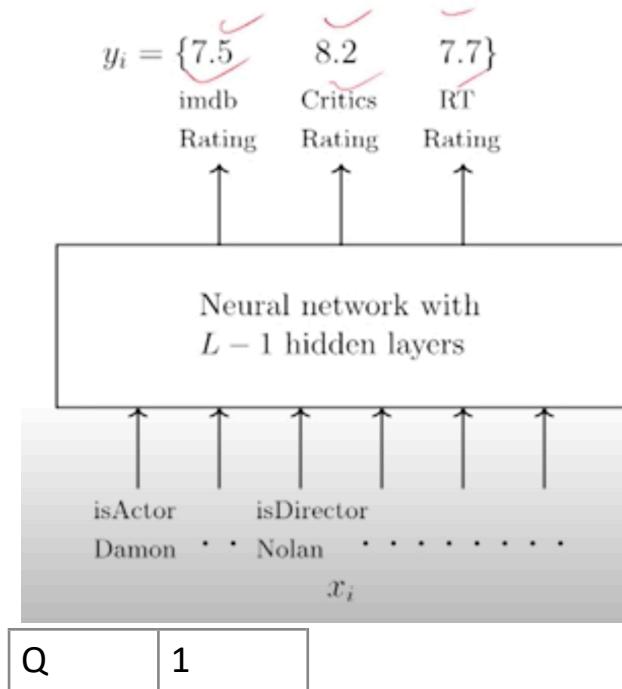


# DL Basics

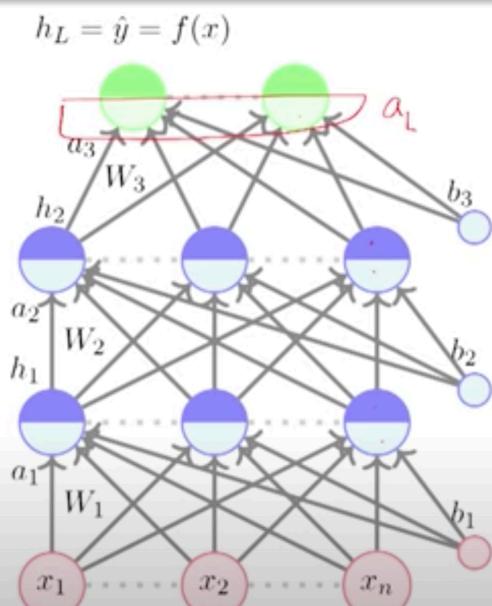
Tuesday, 28 December 2021 2:32 PM

## Output functions and Loss functions in FFNN:



- 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

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



- 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 linear function

$$f(x) = h_L = O(a_L) \\ = W_O a_L + b_O$$



$$y = \begin{bmatrix} 0.35 & 0.25 & 0.4 \end{bmatrix} \quad \checkmark$$

$$\hat{y} = \begin{bmatrix} 0.25 & 0.45 & 0.3 \end{bmatrix} \quad \checkmark$$

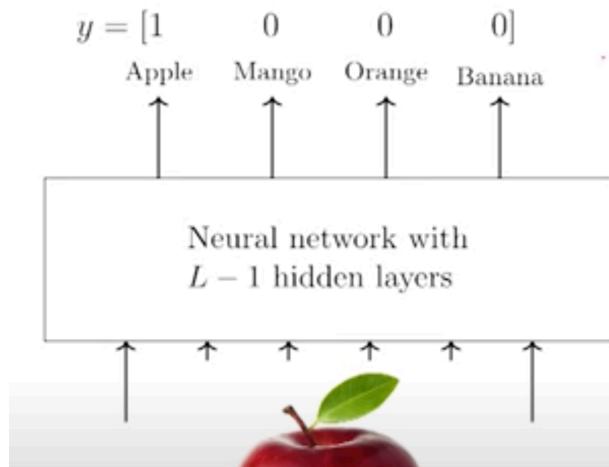
$$\sum_{i=1}^4 q(i) V(i)$$



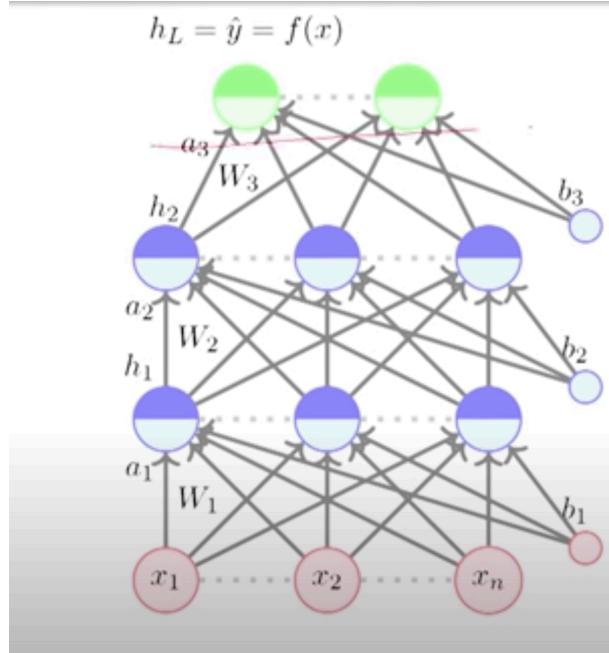
P  
q<sub>i</sub>







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

$$\hat{y}_j = O(a_L)_j = \frac{e^{a_{L,j}}}{\sum_{i=1}^k e^{a_{L,i}}}$$

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

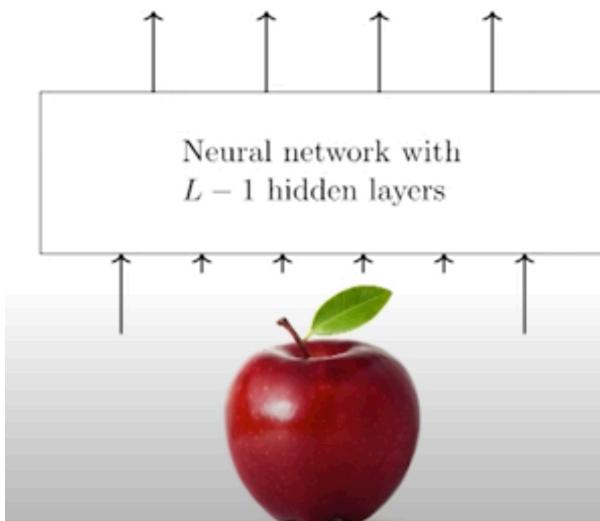
$$a_1 = \begin{bmatrix} 10 & -20 & 30 \\ a_{11} & a_{12} & a_{13} \end{bmatrix} \quad \hat{y}_1 = \frac{e^{10}}{e^{10} + e^{20} + e^{30}}$$



- 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?

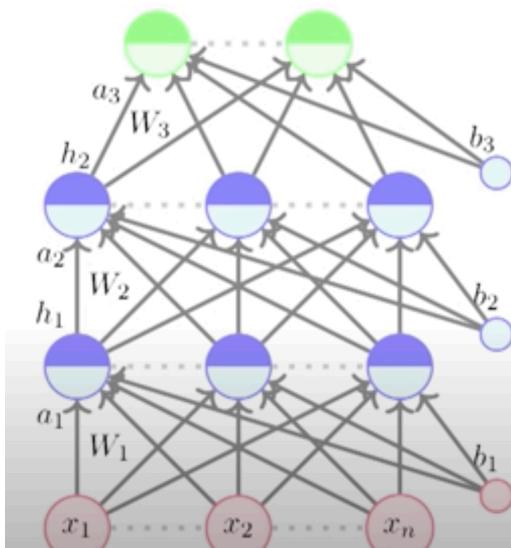
- Cross-entropy

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

- Notice that

$$\begin{aligned} y_c &= 1 && \text{if } c = \ell \text{ (the true class label)} \\ &= 0 && \text{otherwise} \\ \therefore \mathcal{L}(\theta) &= -\log \hat{y}_\ell \end{aligned}$$

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



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

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

- But wait!
- Is  $\hat{y}_\ell$  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  $\theta$
- What does  $\hat{y}_\ell$  encode?
- It is the probability that  $x$  belongs to the  $\ell^{th}$  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	Linear	Softmax
Loss Function	Squared Error	Cross Entropy

$$\frac{\partial \mathcal{L}(\theta)}{\partial \theta}$$

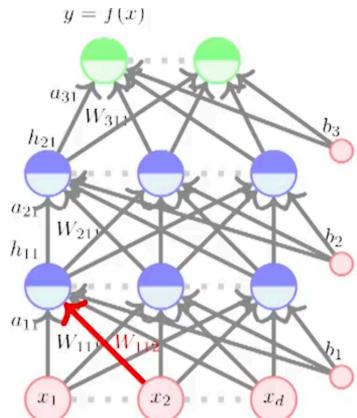
- Of course, there could be other loss functions depending on the problem at hand but the two loss functions that we just saw are encountered very often
- For the rest of this lecture we will focus on the case where the output activation is a softmax function and the loss function is cross entropy

## 4.1) BACK PROPAGATION:





- Let us focus on this one weight ( $W_{112}$ ).
- To learn this weight using SGD we need a formula for  $\frac{\partial \mathcal{L}(\theta)}{\partial W_{112}}$ .



**Algorithm:** gradient descent()

---

```

t ← 0;
max_iterations ← 1000;
Initialize θ₀;
while t++ < max_iterations
do
    θt+1 ← θt - η∇θt;
end

```

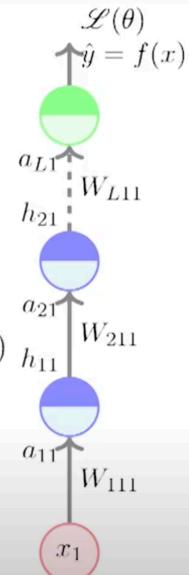
- First let us take the simple case when we have a deep but thin network.
- In this case it is easy to find the derivative by chain rule.

$$\frac{\partial \mathcal{L}(\theta)}{\partial W_{111}} = \frac{\partial \mathcal{L}(\theta)}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial a_{L11}} \frac{\partial a_{L11}}{\partial h_{21}} \frac{\partial h_{21}}{\partial a_{21}} \frac{\partial a_{21}}{\partial h_{11}} \frac{\partial h_{11}}{\partial a_{11}} \frac{\partial a_{11}}{\partial W_{111}}$$

$$\frac{\partial \mathcal{L}(\theta)}{\partial W_{111}} = \frac{\partial \mathcal{L}(\theta)}{\partial h_{11}} \frac{\partial h_{11}}{\partial W_{111}} \quad (\text{just compressing the chain rule})$$

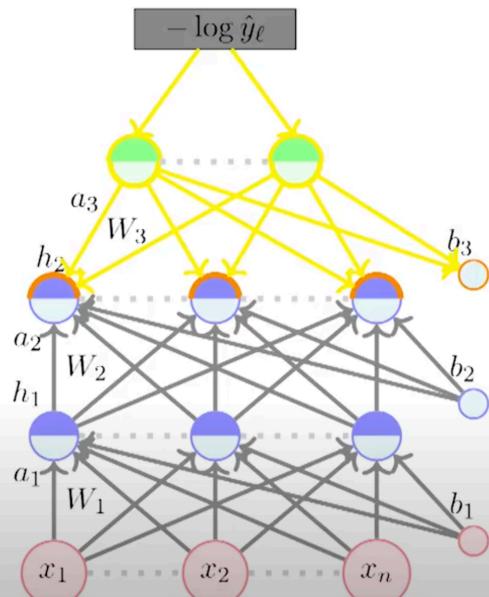
$$\frac{\partial \mathcal{L}(\theta)}{\partial W_{211}} = \frac{\partial \mathcal{L}(\theta)}{\partial h_{21}} \frac{\partial h_{21}}{\partial W_{211}}$$

$$\frac{\partial \mathcal{L}(\theta)}{\partial W_{L11}} = \frac{\partial \mathcal{L}(\theta)}{\partial a_{L1}} \frac{\partial a_{L1}}{\partial W_{L11}}$$



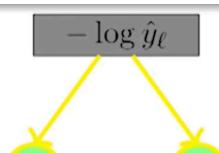
- We get a certain loss at the output and we try to figure out who is responsible for this loss
- So, we talk to the output layer and say "Hey! You are not producing the desired output, better take responsibility".
- The output layer says "Well, I take responsibility for my part but please understand that I am only as good as the hidden layer and weights below me". After all ...

$$f(x) = \hat{y} = O(W_L h_{L-1} + b_L)$$



Let us first consider the partial derivative w.r.t.  $i$ -th output

$$\mathcal{L}(\theta) = -\log \hat{y}_\ell \quad (\ell = \text{true class label})$$



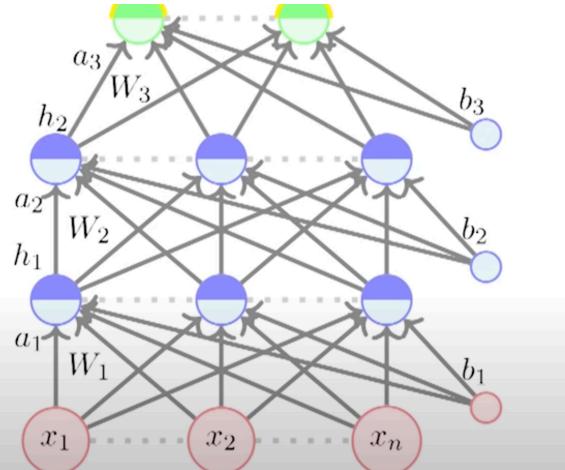




$$\begin{aligned}\frac{\partial}{\partial \hat{y}_i} (\mathcal{L}(\theta)) &= \frac{\partial}{\partial \hat{y}_i} (-\log \hat{y}_\ell) \\ &= -\frac{1}{\hat{y}_\ell} \quad \text{if } i = \ell \\ &= 0 \quad \text{otherwise}\end{aligned}$$

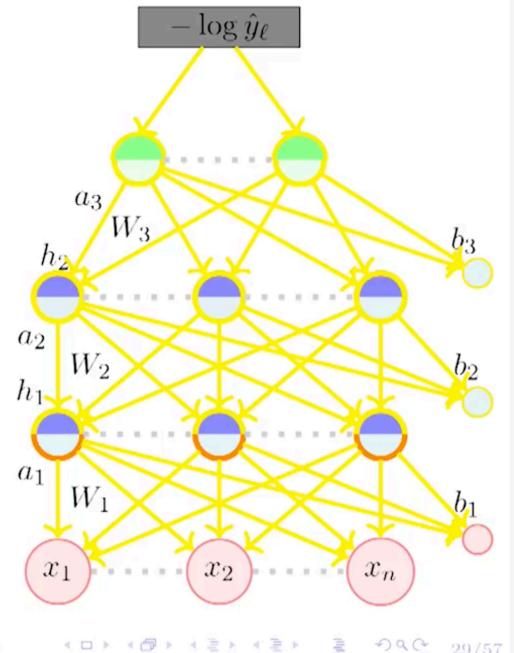
More compactly,

$$\frac{\partial}{\partial \hat{y}_i} (\mathcal{L}(\theta)) = -\frac{\mathbb{1}_{(i=\ell)}}{\hat{y}_\ell}$$



- So, we talk to  $W_L, b_L$  and  $h_L$  and ask them “What is wrong with you?”
- $W_L$  and  $b_L$  take full responsibility but  $h_L$  says “Well, please understand that I am only as good as the pre-activation layer”
- The pre-activation layer in turn says that I am only as good as the hidden layer and weights below me.
- We continue in this manner and realize that the responsibility lies with all the weights and biases (i.e. all the parameters of the model)
- But instead of talking to them directly, it is easier to talk to them through the hidden layers and output layers (and this is exactly what the chain rule allows us to do)

$$\underbrace{\frac{\partial \mathcal{L}(\theta)}{\partial W_{111}}}_{\text{Talk to the weight directly}} = \underbrace{\frac{\partial \mathcal{L}(\theta)}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial a_3}}_{\text{Talk to the output layer}} \underbrace{\frac{\partial a_3}{\partial h_2} \frac{\partial h_2}{\partial a_2}}_{\text{Talk to the previous hidden layer}} \underbrace{\frac{\partial a_2}{\partial h_1} \frac{\partial h_1}{\partial a_1}}_{\text{Talk to the previous hidden layer}} \underbrace{\frac{\partial a_1}{\partial W_{111}}}_{\text{and now talk to the weights}}$$



### Quantities of interest (roadmap for the remaining part):

- Gradient w.r.t. output units
- Gradient w.r.t. hidden units
- Gradient w.r.t. weights and biases

$$\underbrace{\frac{\partial \mathcal{L}(\theta)}{\partial W_{111}}}_{\text{Talk to the weight directly}} = \underbrace{\frac{\partial \mathcal{L}(\theta)}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial a_3}}_{\text{Talk to the output layer}} \underbrace{\frac{\partial a_3}{\partial h_2} \frac{\partial h_2}{\partial a_2}}_{\text{Talk to the previous hidden layer}} \underbrace{\frac{\partial a_2}{\partial h_1} \frac{\partial h_1}{\partial a_1}}_{\text{Talk to the previous hidden layer}} \underbrace{\frac{\partial a_1}{\partial W_{111}}}_{\text{and now talk to the weights}}$$

- Our focus is on *Cross entropy loss* and *Softmax output*.











