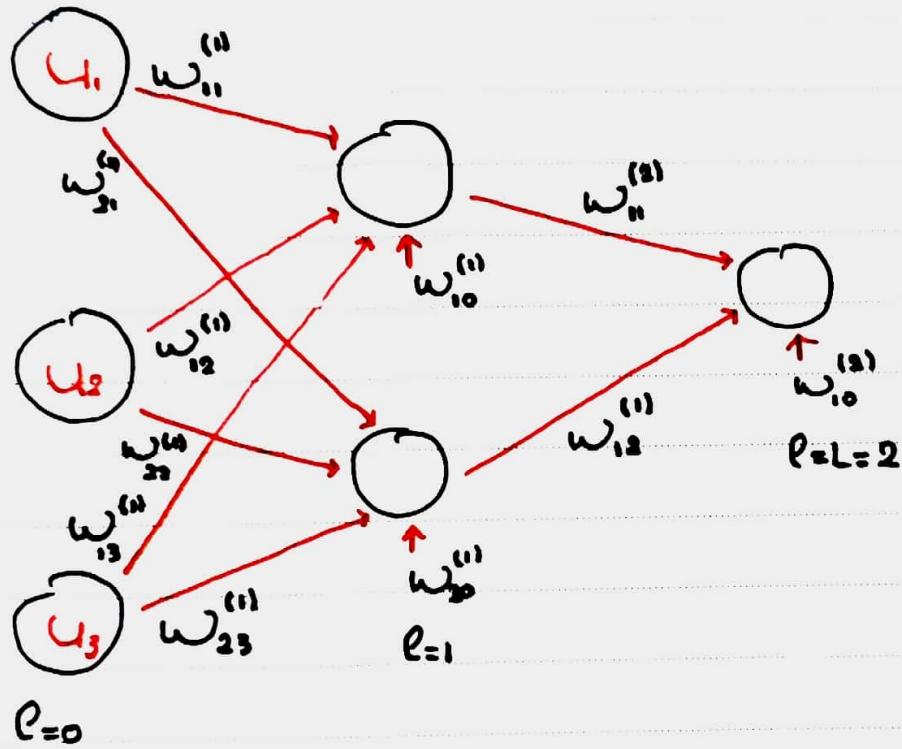


NN Sheet 5

1) Consider

- Assume familiarity with lecture notation
- Drop (in) which denotes current batch



Recall \downarrow Going to (layer)

$w_{ij}^{(l)}$ \nwarrow Coming going from to (neuron)

\rightarrow Can make a $W_{N(l) \times N(l-1)}$ matrix
For each layer $l=1, 2$

let
 $W_0^{(l)} = \begin{pmatrix} w_{10}^{(l)} \\ w_{20}^{(l)} \end{pmatrix}_{N(l)=1}$

- Recall that the i th neuron in the l th layer has output

$$Y_i^{(l)} = f(X_i^{(l)})$$

where

$$X_i^{(l)} = \sum_{j=1}^{N(l)} (w_{ij}^{(l)} \cdot Y_j^{(l-1)}) + w_{i0}^{(l)}$$

Hence, by matrix multiplication if we let $\underline{Y} = \begin{pmatrix} Y_1^{(l)} \\ Y_2^{(l)} \end{pmatrix}$ (and similarly for $\underline{X}^{(l)}$) then

- $\underline{X}^{(l)} = W^{(l)} \underline{Y}^{(l-1)} + W_0^{(l)}$

- $\underline{Y}^{(l)} = f(\underline{X}^{(l)})$

• where f is applied element wise.

q) Equations of Feed Forward Propagation (i.e., X and Y for each layer)

$l=0$)

$$\left. \begin{aligned} Y_1^{(0)} &= U_0 \\ Y_2^{(0)} &= U_1 \\ Y_3^{(0)} &= U_2 \end{aligned} \right\} \quad \underline{Y}^{(0)} = \underline{U}$$

$l=1$)

$$\left. \begin{aligned} X_1^{(1)} &= \sum_{j=1}^3 (\omega_{1j}^{(1)} \cdot U_j) + \omega_{10}^{(1)} \\ Y_1^{(1)} &= P(X_1^{(1)}) \\ X_2^{(1)} &= \sum_{j=1}^3 (\omega_{2j}^{(1)} \cdot U_j) + \omega_{20}^{(1)} \\ Y_2^{(1)} &= P(X_2^{(1)}) \end{aligned} \right\} \quad \begin{aligned} \underline{X}^{(1)} &= W \underline{U} + \underline{w}^{(1)} \\ \underline{Y}^{(1)} &= P(\underline{X}^{(1)}) \end{aligned}$$

$l=2$)

$$\left. \begin{aligned} X_1^{(2)} &= \sum_{j=1}^2 (\omega_{1j}^{(2)} \cdot Y_j^{(2)}) + \omega_{10}^{(2)} \\ Y_1^{(2)} &= P(X_1^{(2)}) \end{aligned} \right\} \quad \begin{aligned} \underline{X}^{(2)} &= W \underline{Y}^{(2)} + \underline{w}^{(2)} \\ \underline{Y}^{(2)} &= P(\underline{X}^{(2)}) \end{aligned}$$

- On the left are the equations in index notation (by the definition of the neuron) and on the right are the equivalent vector representation for the whole layer (as was discussed last Page).

- Note that the activation function is given as

$$P(x) = \frac{1}{1+e^{-x}} \quad (\text{Sigmoid Function})$$

- b) Derive the loss for binary classification

- The neural network's output model is

$$P(y=1|\theta)$$

↑ network parameters

- Hence, y follows the Bernoulli distribution

$$P(y_i|\theta) = (y_i P(y_i|\theta))^y (1 - P(y_i|\theta))^{1-y}$$

when $y=1$ when $y=0$

- The Probability of observing the dataset is

$$\prod_{i=1}^m P(y_i|\theta)$$

Maximizing this is equal to maximizing
 $\sum_{i=1}^m \log P(y_i|\theta)$

- To find θ that maximizes such probability, maximize

(

$$\sum_{i=1}^m \ln P(y_i|\theta) = \sum_{i=1}^m y_i \ln(P(y_i=1|\theta)) + (1-y_i) \ln(1-P(y_i=1|\theta))$$

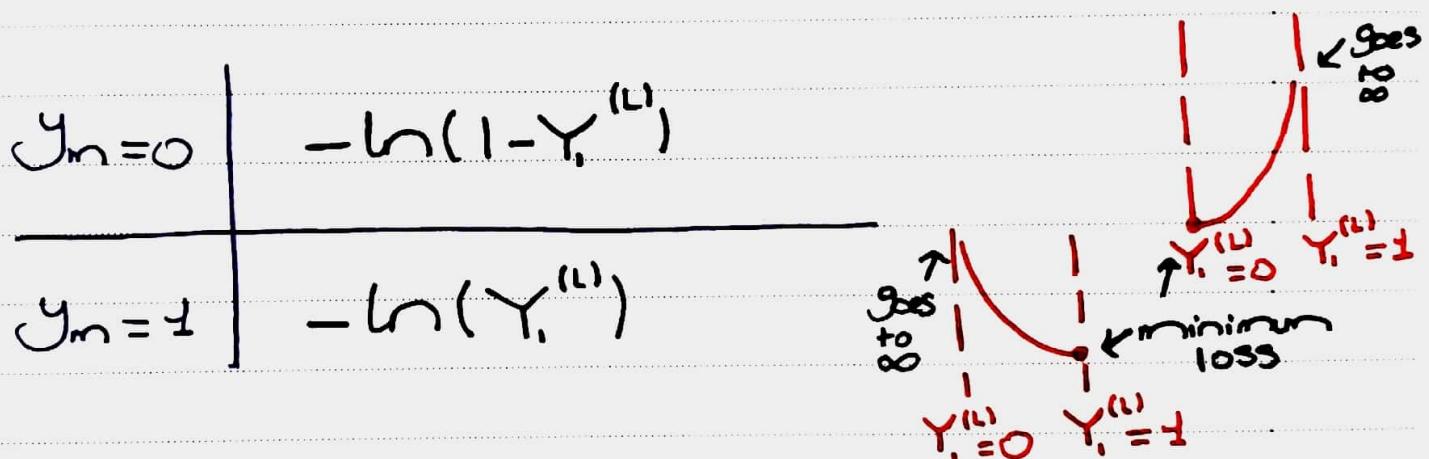
which is equivalent to minimizing

$$J = -\frac{1}{M} \sum_{i=1}^M y_i \ln(P(y_i=1|\theta)) + (1-y_i) \ln(1-P(y_i=1|\theta))$$

To yield the Parameters θ for which the Observed data is most Probable.

- Known as binary Cross entropy and can be justified using two ways

$$1-\text{Consider } J_m = -y_m \ln Y_i^{(L)} - (1-y_m) \ln(1-Y_i^{(L)})$$



- The loss is hence 0 for correct classification and gets higher the more incorrect it is

q , relative to P

- 2- The Cross entropy between 2 Probability distributions P and q is given by

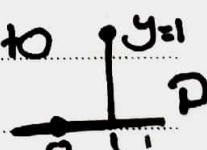
$$H(P, q) = - \sum_{x \in X} P(x) \log(q(x))$$

• gives entropy when $q(x) = P(x)$

- Measures the Similarity between two Probability distributions (becomes entropy (min.) when they're the same)

For each example
• Minimise total cross entropy

Here we want $Y^{(L)}$ to be as close as possible to q .



c) Derive the equations of backpropagation

$$\boxed{J_m}$$

Recall,

$$J = \underbrace{\sum_{m=1}^M -y_m \ln Y_i^{(1)} - (1-y_m) \ln (1-Y_i^{(1)})}_{J_m}$$

$$X_i^{(2)} = \sum_{j=1}^2 (w_{ij}^{(2)} \cdot Y_j^{(1)}) + w_{i0}^{(2)} \quad \left. \begin{array}{l} \text{last layer} \\ l=2=L \end{array} \right\}$$

$$Y_i^{(2)} = P(X_i^{(2)})$$

$$X_i^{(1)} = \sum_{j=1}^3 (w_{ij}^{(1)} \cdot Y_j^{(0)}) + w_{i0}^{(1)} \quad \left. \begin{array}{l} l=1 \\ i=1,2 \end{array} \right\}$$

$$Y_i^{(1)} = P(X_i^{(1)})$$

$$Y_i^{(0)} = U_i$$

$$\left. \begin{array}{l} l=0 \\ i=1,2,3 \end{array} \right\}$$

→ Let's Find $\frac{\partial J_m}{\partial w_{ij}^{(2)}}$ • $\frac{\partial J}{\partial w} = \sum_{m=1}^M \frac{\partial J_m}{\partial w}$

$$\cdot \frac{\partial J_m}{\partial w_{ij}^{(2)}} = \frac{\partial J_m}{\partial Y_i^{(2)}} \cdot \frac{\partial Y_i^{(2)}}{\partial X_i^{(2)}} \cdot \frac{\partial X_i^{(2)}}{\partial w_{ij}^{(2)}}$$

$$\textcircled{1} \frac{\partial J_m}{\partial Y_i^{(2)}} = -y_m \cdot \frac{1}{Y_i^{(2)}} - (1-y_m) \cdot \frac{-1}{1-Y_i^{(2)}} = \begin{cases} \frac{1}{1-Y_i^{(2)}} & y_m=0 \\ \frac{-1}{Y_i^{(2)}} & y_m=1 \end{cases}$$

$$\textcircled{2} \frac{\partial Y_i^{(2)}}{\partial X_i^{(2)}} = P'(X_i^{(2)})$$

$$\text{Since } P(x) = \frac{1}{1+e^{-x}} \text{ then } P'(x) = -(1+e^{-x})^{-2} \cdot e^{-x}$$

$$= \frac{e^{-x}}{1+e^{-x}} \cdot \frac{1}{1+e^{-x}}$$

$$= \frac{1}{1+e^{-x}} \cdot \left(\frac{1+e^{-x}}{1+e^{-x}} - \frac{1}{1+e^{-x}} \right)$$

$$= P(x)(1-P(x))$$

Thus,

$$\frac{\partial Y_i^{(2)}}{\partial X_i^{(2)}} = Y_i^{(2)}(1-Y_i^{(2)})$$

$$③ \frac{\partial X_i^{(2)}}{\partial w_{ij}^{(2)}} = \frac{\partial}{\partial w_{ij}^{(2)}} \left(\sum_{j=1}^2 (w_{ij}^{(2)} Y_j^{(1)}) + b_j^{(2)} \right)$$

$$= \frac{\partial}{\partial w_{ij}^{(2)}} (w_{ij}^{(2)} Y_j^{(1)}) = Y_j^{(1)} \quad j=1,2 \\ \cdot Y_0^{(1)} = 1$$

Thus,

$$\frac{\partial J_m}{\partial w_{ij}^{(2)}} = (-y_m + (1-y_m)) \cdot Y_i^{(2)}(1-Y_i^{(2)}) \cdot Y_j^{(1)}$$

$$= (-y_m(1-Y_i^{(2)}) + (1-y_m)Y_i^{(2)}) \cdot Y_j^{(1)}$$

$$= \begin{cases} Y_i^{(2)} \cdot Y_j^{(1)} & y_m=0 \\ -Y_j^{(1)} \cdot (1-Y_i^{(2)}) & y_m=1 \end{cases}$$

I only later noticed
that we can just
write this as
 $Y_j^{(1)}(Y_i^{(2)} - y_m)$

exclude for $j=0$
(bias)

In vector form:

$$\frac{\partial J_m}{\partial W^{(2)}} = \left(\frac{\partial J_m}{\partial w_{11}^{(2)}} \frac{\partial J_m}{\partial w_{12}^{(2)}} \right) = \underline{Y}^{(1)} \cdot \begin{cases} Y_i^{(2)} & y_m=0 \\ -(1-Y_i^{(2)}) & y_m=1 \end{cases}$$

→ Now Find $\frac{\partial J_m}{\partial w_{IJ}^{(1)}}$



$$\frac{\partial J_m}{\partial w_{IJ}^{(1)}} = \underbrace{\frac{\partial J_m}{\partial Y_I^{(2)}}}_{-y_m(1-Y_I^{(2)}) + (1-y_m)Y_I^{(2)}} \cdot \underbrace{\frac{\partial Y_I^{(2)}}{\partial X_I^{(2)}}}_{\text{(calculated Previously)}} \cdot \underbrace{\frac{\partial X_I^{(2)}}{\partial Y_I^{(1)}}}_{w_{I1}^{(2)}} \cdot \underbrace{\frac{\partial Y_I^{(1)}}{\partial X_I^{(1)}}}_{w_{I1}^{(1)}}$$

$$\frac{\partial X_I^{(2)}}{\partial Y_I^{(1)}} = \frac{\partial}{\partial Y_I^{(1)}} \sum_{j=1}^2 (w_{Ij} \cdot Y_j^{(1)}) + w_{I0}^{(2)}$$

$$= \frac{\partial}{\partial Y_I^{(1)}} (w_{I1}^{(2)} \cdot Y_I^{(1)} + w_{I0}^{(2)}) = w_{I1}^{(2)}$$

• regardless
of $y_I^{(1)}$ (bias)

$$\frac{\partial Y_I^{(1)}}{\partial X_I^{(1)}} = Y_I^{(1)}(1-Y_I^{(1)})$$

$$\frac{\partial X_I^{(1)}}{\partial w_{IJ}^{(1)}} = \frac{\partial}{\partial w_{IJ}^{(1)}} \sum_{j=1}^3 (w_{Ij}^{(1)} \cdot u_j) + w_{I0}^{(1)} = \frac{\partial}{\partial w_{IJ}^{(1)}} w_{IJ}^{(1)} \cdot u_j$$

$$\frac{\partial J_m}{\partial w_{IJ}^{(1)}} = \begin{cases} Y_I^{(2)} \cdot w_{I1}^{(2)} \cdot Y_I^{(1)}(1-Y_I^{(1)}) \cdot u_j & = u_j \quad y_m = 0 \\ (Y_I^{(2)} - 1) \cdot w_{I1}^{(2)} \cdot Y_I^{(1)}(1-Y_I^{(1)}) \cdot u_j & = y_m = -1 \end{cases}$$

→ which if we let \odot denote element-wise vector multiplication gives in vector notation:

$$\frac{\partial J_m}{\partial w^{(1)}} = \begin{cases} Y_i^{(2)} \cdot ((w^{(2)})^t \odot Y^{(1)} \odot (1 - Y^{(1)})) \underline{U}^t & y_m=0 \\ (Y_i^{(2)} - 1) \cdot ((w^{(2)})^t \odot Y^{(1)} \odot (1 - Y^{(1)})) \underline{U}^t & y_m=1 \end{cases}$$

We earlier had

$$\frac{\partial J_m}{\partial w^{(2)}} = \begin{cases} Y_i^{(2)} \cdot (Y^{(1)})^t & y_m=0 \\ (Y_i - 1)^{(2)} \cdot (Y^{(1)})^t & y_m=1 \end{cases}$$

• ignore last term for bias vector

→ These are the equations of backward propagation of the network.

- Top of 1st of last 3 Pages includes the forward propagation equations.

d) Suppose an input $\underline{U} = \begin{pmatrix} 3 \\ 0.2 \\ -0.1 \end{pmatrix}$ what's the output

• $l=0$ $\underline{Y}^{(0)} = \underline{U} = \begin{pmatrix} 3 \\ 0.2 \\ -0.1 \end{pmatrix}$

• $l=1$ $\underline{X}^{(1)} = \underline{W}^{(1)} \underline{U} + \underline{W}_0^{(1)} = \begin{pmatrix} 0.31 \\ -0.28 \end{pmatrix}$

$$\begin{pmatrix} 0.1 & 0.2 & 0.3 \\ -0.1 & 0 & -0.2 \end{pmatrix}$$

$$\underline{Y}^{(1)} = F(\underline{X}^{(1)}) = \begin{pmatrix} 0.5769 \\ 0.43045 \end{pmatrix}$$

$$4/(1+e^{-x})^{-1}$$

$\ell = 2$

$$\underline{X}^{(2)} = W^{(2)} \underline{Y}^{(1)} + w_0^{(2)} = 0.02929$$

\downarrow
 $(0.2 - 0.2) \begin{pmatrix} 0.5769 \\ 0.43045 \end{pmatrix}$

$$\underline{Y}^{(2)} = P(\underline{X}^{(2)}) = 0.50732 \quad \leftarrow \text{Network's OutPut.}$$

e) $0.50732 < 0.6$ and Hence, class 0 will be chosen.

P) $J_m = -y_m \ln Y^{(2)} - (1-y_m) \ln (1-Y^{(2)})$
 $\cdot y_m = 0$
 $= -\ln (1-Y^{(2)}) = -\ln (1-0.50732)$
 $= 0.7078954$

g) $\frac{\partial J_m}{\partial W^{(2)}} = Y_i^{(2)} (\underline{Y}^{(1)})^t = 0.50732 (0.5769 \ 0.43045)$
 $= (0.29267 \ 0.218376)$

$$\frac{\partial J_m}{\partial w_0^{(2)}} = Y_i^{(2)} = 0.50732$$

$$\frac{\partial J_m}{\partial W^{(1)}} = Y_i^{(2)} ((W^{(2)})^t \underbrace{\odot Y_0^{(1)} (1-Y_0^{(1)})}_{= 0.2441 / 0.24516}) \underline{Y}^t$$
 $= 0.50732 \left(\begin{pmatrix} 0.2 \\ -0.2 \end{pmatrix} \odot \begin{pmatrix} 0.2441 \\ 0.24516 \end{pmatrix} \right) (3 \ 0.2 - 0.1)$
 $= \begin{pmatrix} 0.02476736 \\ -0.0248749 \end{pmatrix} (3 \ 0.2 - 0.1)$

$$= \begin{pmatrix} 0.0743 & 4.953 \times 10^{-3} & -2.4767 \times 10^{-3} \\ -0.0746 & -4.974 \times 10^{-3} & 2.48749 \times 10^{-3} \end{pmatrix}$$

$$\frac{\partial J_n}{\partial w_0^{(1)}} = Y_1^{(2)} ((W^{(2)})^t \odot Y_0^{(1)} (I - Y_0^{(1)})) = \begin{pmatrix} 0.024767 \\ -0.0248749 \end{pmatrix}$$

. computed last page

h)

$$W_{(\text{new})}^{(2)} = W_{(\text{old})}^{(2)} - \gamma \frac{\partial J_n}{\partial w_{(\text{old})}^{(2)}} \quad \begin{matrix} 0.05 \\ \downarrow \\ (0.2 \quad -0.2) \end{matrix} \rightarrow (0.29267 \quad 0.218376)$$

$$= (0.1853665 \quad -0.2109188)$$

$$W_{(\text{new})}^{(2)} = W_{(\text{old})}^{(2)} - \gamma \frac{\partial J_n}{\partial w_{(\text{old})}^{(2)}} = -0.025366 \quad \begin{matrix} \downarrow \\ 0 \end{matrix} \quad \begin{matrix} \downarrow \\ 0.03 \end{matrix} \quad \rightarrow 0.50732$$

$$W_{(\text{new})}^{(4)} = W_{(\text{old})}^{(4)} - \gamma \frac{\partial J_n}{\partial w_{(\text{old})}^{(4)}} \rightarrow \begin{pmatrix} 0.0743 & 4.953 \times 10^{-3} & -2.4767 \times 10^{-3} \\ -0.0746 & -4.974 \times 10^{-3} & 2.48749 \times 10^{-3} \end{pmatrix}$$

$$(0.1 \quad 0.2 \quad 0.3) \\ (-0.1 \quad 0 \quad -0.2)$$

$$= \begin{pmatrix} 0.096285 & 0.199752 & 0.30124 \\ -0.09627 & 2.487 \times 10^{-4} & -0.8 \end{pmatrix}$$

$$W_{(new)}^{(2)} = W_{(old)}^{(2)} - \eta \cdot \frac{\partial J_m}{\partial W_{(old)}^{(2)}}$$

↓ ↓
 $\begin{pmatrix} 0 \\ 0 \end{pmatrix}$ 0.05 $\rightarrow \begin{pmatrix} 0.024767 \\ -0.0248749 \end{pmatrix}$

$$= \begin{pmatrix} -1.23835 \times 10^{-3} \\ -1.243745 \times 10^{-3} \end{pmatrix}$$

2) If it was solving a regression problem where the output can take any value then the activation function at the output layer must be linear i.e. $P(x) = x, P'(x) = 1$

• Feedforward Equations

$$\underline{Y}^{(0)} = \underline{U}$$

$$\underline{X}^{(1)} = W^{(1)} \underline{U} + W_0^{(1)}$$

$$\underline{Y}^{(1)} = P(\underline{X}^{(1)})$$

. Hidden layer activations
need not change

$$\underline{X}^{(2)} = W^{(2)} \underline{Y}^{(1)} + W_0^{(2)}$$

$$\underline{Y}^{(2)} = \underline{X}^{(2)} \quad \leftarrow \text{Only Change in FeedForward}$$

* The loss function should change as well to a regression loss:

$$J_m = \frac{1}{2} (y_m - Y^{(2)})^2, \quad J = \frac{1}{M} \sum_{m=1}^M J_m \quad (\text{MSE})$$

• Backward Propagation Equations

$$\rightarrow \frac{\partial J_m}{\partial w_{ij}^{(2)}} = \frac{\partial J_m}{\partial Y_i^{(2)}} \cdot \frac{\partial Y_i^{(2)}}{\partial X_i^{(1)}} \cdot \frac{\partial X_i^{(2)}}{\partial w_{ij}^{(1)}}$$

$(Y^{(2)} - y_m)$ $\downarrow Y_i$ $\downarrow Y^{(1)}$

→ if you look closely
 $\frac{\partial J_m}{\partial Y_i^{(2)}} \cdot \frac{\partial Y_i^{(2)}}{\partial X_i^{(1)}}$
 didn't really change

$$= (Y^{(2)} - y_m) \underline{Y^{(1)}} \quad \text{↓ exclude for bias (w.)}$$

$$\frac{\partial J_m}{\partial w^{(2)}} = (Y^{(2)} - y_m) (\underline{Y^{(1)}})^t$$

$$\rightarrow \frac{\partial J_m}{\partial w_{ij}^{(1)}} = \underbrace{\frac{\partial J_m}{\partial Y_i^{(2)}} \cdot \frac{\partial Y_i^{(2)}}{\partial X_i^{(2)}} \cdot \frac{\partial X_i^{(2)}}{\partial Y_I^{(1)}} \cdot \frac{\partial Y_I^{(1)}}{\partial X_I^{(1)}} \cdot \frac{\partial X_I^{(1)}}{\partial w_{ij}^{(1)}}}_{(Y^{(2)} - y_m)} \cdot \underline{w_{jI}^{(2)}} \cdot \underline{Y_I^{(1)}(1 - Y_I^{(1)})} \cdot \underline{u_j}$$

$$= (Y^{(2)} - y_m) ((w^{(2)})^t \odot \underline{Y} \odot \underline{(1 - Y)} \odot \underline{u})^t \quad \text{↓ exclude for bias}$$

* So in essence, the backpropagation equations won't change!

+ Sigmoid	Linear
+ Binary cross entropy loss	↔ + MSE loss

3) what would be changed for multi-class classification?

- The Output layer should involve K neurons rather than 1 (one for each class)
- The labels in the dataset are now Kx1 vectors if the ith index = 1 then the example belongs to ith class

- The activations used in output layer should be Softmax

$$Y_i^{(L)} = \frac{e^{x_i^{(L)}}}{\sum_{j=1}^K e^{x_j^{(L)}}}$$

• Looks like we'll have to revisit chain rule



- The loss function should generalize from binary cross entropy to multi-class cross entropy

$$J_m = - \sum_{i=1}^K y_{m,i} \ln(Y_i^{(2)})$$

* Changes in FeedForward Pass

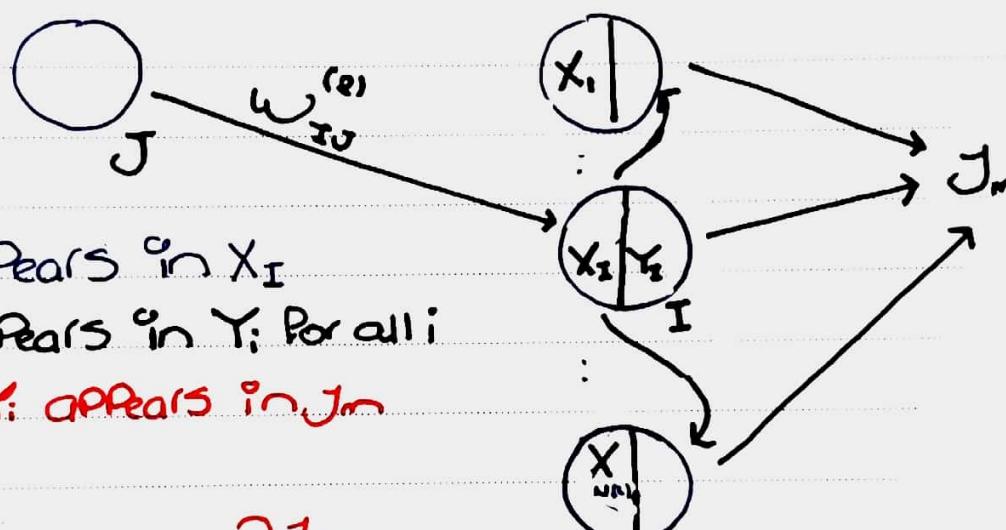
$$X^{(2)} = W_{K \times 2}^{(2)} Y^{(1)} + W_{0 \times K \times 1}^{(2)} \quad (X_i^{(2)} = \sum_{j=1}^2 w_{ij}^{(2)} Y_j^{(1)} + w_{i0}^{(2)})$$

$$Y^{(2)} = e^{X^{(2)}} \cdot \frac{1}{\sum_{j=1}^K e^{x_j^{(2)}}} \quad (Y_i^{(2)} = \frac{e^{x_i^{(2)}}}{\sum_{j=1}^K e^{x_j^{(2)}}})$$

* Changes in backward Pass

→ First want $\frac{\partial J_m}{\partial w_{IJ}^{(2)}}$

Observe



- $w_{IJ}^{(2)}$ appears in X_I
- which appears in Y_i for all i
- and each Y_i appears in J_m

$$\underbrace{\frac{\partial J_m}{\partial X_I^{(2)}}}_{\text{Red}}$$

$$\frac{\partial J_m}{\partial w_{IJ}^{(2)}} = \left(\sum_{i=1}^k \frac{\partial J_m}{\partial Y_i^{(2)}} \cdot \frac{\partial Y_i^{(2)}}{\partial X_I^{(2)}} \right) \cdot \frac{\partial X_I}{\partial w_{IJ}^{(2)}}$$

$$\begin{aligned} \frac{\partial}{\partial w_{IJ}^{(2)}} \sum_{j=1}^J w_{IJ}^{(2)} Y_j^{(1)} + w_{IJ}^{(2)} \\ = Y_J^{(1)} \quad (\text{only for } J \neq 0) \end{aligned}$$

$$\underbrace{\frac{\partial J_m}{\partial X_I^{(2)}}}_{\text{Red}}$$

$$\sum_{i=1}^k \frac{\partial J_m}{\partial Y_i^{(2)}} \cdot \frac{\partial Y_i^{(2)}}{\partial X_I^{(2)}} + \frac{\partial J_m}{\partial Y_I^{(2)}} \cdot \frac{\partial Y_I^{(2)}}{\partial X_I^{(2)}}$$

$$① \quad \frac{\partial J_m}{\partial Y_I^{(2)}} = \frac{\partial}{\partial Y_I^{(2)}} \left(-\sum_{i=1}^k y_{m,i} \ln(Y_i^{(2)}) \right) = \frac{\partial}{\partial Y_I^{(2)}} - y_{m,I} \ln(Y_I^{(2)}) = -\frac{y_{m,I}}{Y_I^{(2)}}$$

$$\begin{aligned}
 ② \frac{\partial Y_I^{(2)}}{\partial X_I^{(2)}} &= \frac{\partial}{\partial X_I^{(2)}} \left(\frac{e^{X_I^{(2)}}}{\sum_{i=1}^k e^{X_i^{(2)}}} \right) \\
 &= \frac{e^{X_I^{(2)}}}{\sum_{i=1}^k e^{X_i^{(2)}}} \left(\frac{e^{X_I^{(2)}}}{\sum_{i=1}^k e^{X_i^{(2)}}} - \frac{e^{X_I^{(2)}}}{\sum_{i=1}^k e^{X_i^{(2)}}} \right) \\
 &= Y_I^{(2)}(1 - Y_I^{(2)})
 \end{aligned}$$

• Take log
then differentiate

$$③ \frac{\partial J_m}{\partial Y_i^{(2)}} = -y_{m,i} \frac{1}{Y_i^{(2)}}$$

$$\begin{aligned}
 ④ \frac{\partial Y_i^{(2)}}{\partial X_I^{(2)}} &= \frac{\partial}{\partial X_I^{(2)}} \left(\frac{e^{X_i^{(2)}}}{\sum_{j=1}^k e^{X_j^{(2)}}} \right) = \frac{e^{X_i^{(2)}}}{\sum_{j=1}^k e^{X_j^{(2)}}} \left(0 - \frac{e^{X_I^{(2)}}}{\sum_{j=1}^k e^{X_j^{(2)}}} \right) \\
 &= -Y_i^{(2)} \cdot Y_I^{(2)}
 \end{aligned}$$

Now

$$\begin{aligned}
 \frac{\partial J_m}{\partial w_{IJ}^{(2)}} &= \sum_{i=1}^k \frac{\partial J_m}{\partial Y_i^{(2)}} \cdot \frac{\partial Y_i^{(2)}}{\partial X_I^{(2)}} + \frac{\partial J_m}{\partial Y_I^{(2)}} \cdot \frac{\partial Y_I^{(2)}}{\partial X_I^{(2)}} \\
 &= \sum_{i=1}^k -y_{m,i} \cdot \frac{-Y_i^{(2)} \cdot Y_I^{(2)}}{Y_i^{(2)}} + \frac{-y_{m,I}}{Y_I^{(2)}} \cdot Y_I^{(2)}(1 - Y_I^{(2)}) \\
 &= Y_I^{(2)} \sum_{i=1}^k y_{m,i} - y_{m,I} (Y_I^{(2)} - 1) \\
 &= Y_I^{(2)} \sum_{i=1}^k y_{m,i} - y_{m,I} = Y_I^{(2)} - y_{m,I}
 \end{aligned}$$

-1 as y_m is one-hot vector

Thus we have

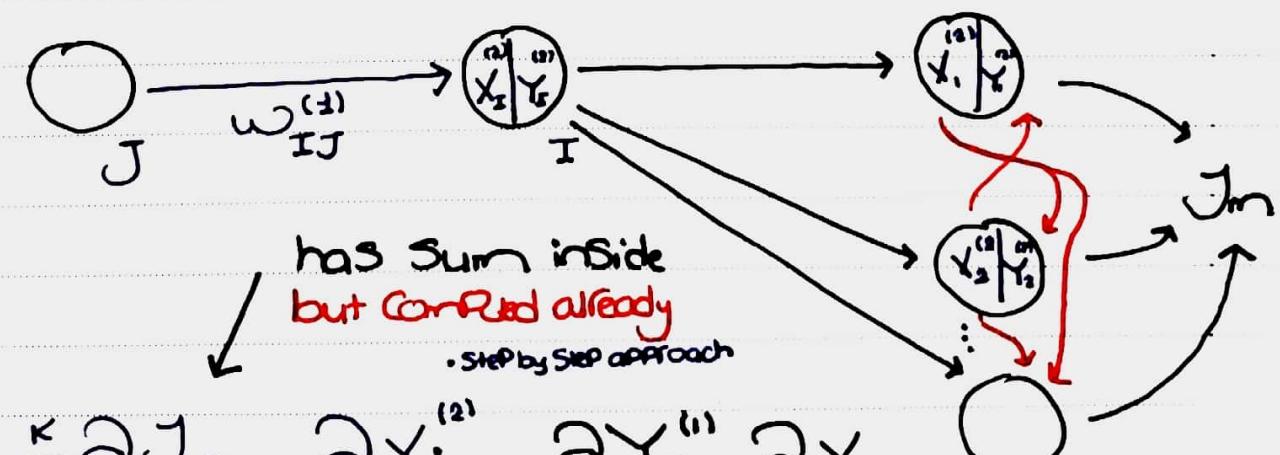
$$\frac{\partial J_m}{\partial w_{IJ}^{(2)}} = (Y_I^{(2)} - y_{m,I}) \cdot Y_J^{(1)}$$

• last term excluded
for bias.

$$\frac{\partial J_m}{\partial w^{(2)}} = (\underline{Y}^{(2)} - \underline{y}_m) (\underline{Y}^{(1)})^t$$

* which is pretty
much the SAME
EXPRESSION
that we had
before (with the
1st term being
generalized to a
vector)

Now For $\frac{\partial J_m}{\partial w_{IJ}^{(1)}}$



$$\begin{aligned} \frac{\partial J_m}{\partial w_{IJ}^{(2)}} &= \sum_{i=1}^k \frac{\partial J_m}{\partial x_i^{(2)}} \cdot \frac{\partial x_i^{(2)}}{\partial y_I^{(1)}} \cdot \frac{\partial y_I^{(1)}}{\partial x_I^{(1)}} \cdot \frac{\partial x_I^{(1)}}{\partial w_{IJ}^{(1)}} \rightarrow u_J \\ &= (Y_I^{(2)} - y_{m,i}) \frac{\partial}{\partial y_I^{(1)}} \left(\sum_{j=1}^2 w_{ij}^{(2)} y_j^{(1)} + w_{i0}^{(2)} \right) = w_{iI}^{(2)} \\ &= \sum_{i=1}^k (Y_I^{(2)} - y_{m,i}) w_{iI}^{(2)} \cdot Y_I^{(1)} (1 - Y_I^{(1)}). u_J \end{aligned}$$

} just like before

exclude for bias

In vector notation

let's call this $\delta_i^{(2)}$

$$\frac{\partial J_m}{\partial W_{IJ}^{(2)}} = \sum_{i=1}^K (\underline{Y}_i^{(2)} - \underline{y}_{m,i}) \underline{w}_{i,I}^{(2)} \cdot (\underline{Y}_I^{(1)} (1 - \underline{Y}_I^{(1)})) \cdot \underline{u}_J$$

- Dot Product between Ith column in weight matrix and $(\underline{Y}^{(2)} - \underline{y}_m)$ vector
- Transmogrify by Ith element here

Now by definition of outer product ($A_{ij} = a_i b_j = \underline{a} \underline{b}^t$)

$$\rightarrow \frac{\partial J_m}{\partial W_{IJ}^{(2)}} = \underline{\delta}^{(2)} \cdot \underline{u}_I^t$$

$$\underline{\delta}^{(2)} = ((\underline{Y}^{(2)} - \underline{y}_m) \underline{W}_{K \times 2}^{(2)})^t \quad \text{need it} \quad \text{a col. vector}$$

- For the Ith element in $\underline{\delta}$, we dot product $(\underline{Y}^{(2)} - \underline{y}_m)$ with the Ith column in \underline{W}

* Observe

If \underline{c} is a $l \times 1$ vector and \underline{a} is a $K \times 1$ vector and B is a $K \times l$ matrix

• Let $C_i = \sum_{j=1}^k a_j B_{ji} = \underline{a}^T B_{:,i} = B_{:,i} \underline{a}$

then $\underline{C} = B^T \underline{a}^T = \begin{pmatrix} \underline{C}_1 \\ \vdots \\ \underline{C}_l \end{pmatrix}$

• Back Prop Eqns:

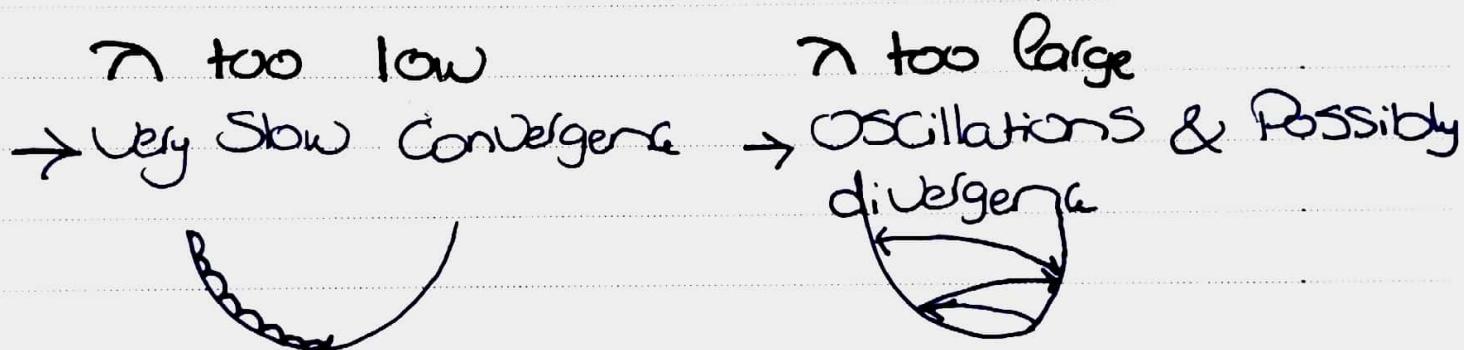
$$\frac{\partial J_m}{\partial W^{(2)}} = (\underline{Y}^{(2)} - \underline{y}_m) (\underline{Y}^{(1)})^t$$

$$\frac{\partial J_m}{\partial W^{(1)}} = ((W^{(2)})^T (\underline{Y}^{(2)} - \underline{y}_m) \odot \underline{Y}^{(1)} (1 - \underline{Y}^{(1)})) \cdot \underline{u}^t$$

• avoid bias term for bias

4) The learning rate signifies the step size used in each iteration of gradient descent to reach the minimum

- Changing η affects convergence in the training process



- Hence, the larger η is the faster is possibly the convergence but the more prone we are to overshooting/divergence.

5) Mini-batch gradient descent divides the training data into a set of batches and performs a weight update batch by batch (the loss function considers one batch at a time).

Mini-batch GD with

Stochastic Gradient Descent \rightarrow Set Batch Size = 1

Full Batch Gradient Descent \rightarrow Set Batch Size = N

Pastel Convergence & less memory Foot Print among other desirable properties.

6) Mention Four types of activation Function

1-

$$P(x) = \text{ReLU}(x) \\ = \max(0, x)$$

$$P'(x) = \begin{cases} 1 & x > 0 \\ 0 & x < 0 \end{cases}$$

→ Simple and easily computed (both $P(x)$ & $P'(x)$)

→ Can cause the "dying ReLU" Problem

- If the inputs to a neuron are -ve for all training data then the derivative will be zero & the weights will never be updated

* leaky ReLU

is a solution.

$$P(x) = \begin{cases} x & x > 0 \\ ax & x < 0 \end{cases}$$

→ Practical results show fast learning (due to const. P') and a Polar Choke when dealing with images.

2-

$$P(x) = x$$

. Linear

activation

$$P'(x) = 1$$

→ Used in the output layer in regression tasks (Predict continuous quantity)

→ Simple & efficient but the neural network becomes linear if used in hidden layers (not a universal approximator).

3-

$$P(x) = \frac{1}{1 + e^{-x}}$$

Sigmoid

$$\dot{P}(x) = P(x)(1 - P(x))$$

• as demonstrated
earlier

- Used in the Output layer when the classification task is binary

- Slow updates (Weak gradients)
- Suffers from the vanishing gradient problem
- Exponential operations

4-

$$P(x) = \tanh(x)$$

$$\dot{P}(x) = \text{Sech}^2(x)$$

• Pretty much Scaled Sigmoid ($2\sigma(2x) - 1$) such that the gradients are stronger (Faster Convergence) but some disadvantages.

5-

$$P(x_i) = \frac{e^{x_i}}{\sum_{j=1}^{n_{\text{cls}}} e^{x_j}}$$

$$\dot{P}_{x_i}(x_i) = P(x_i)(1 - P(x_i))$$

• as demonstrated
earlier

$$\dot{P}_{x_i}(x_i) = -P(x_i) \cdot P(x_i)$$

- Used in the Output layer when the task is multi-class classification.
- Can cause numerical instability, not so efficient...
→ Solved by log Softmax

7) Gradient Descent Hyperparameters

- Depending on the iterative scheme the hyperparameters can range from just the learning rate (and batch size, #epochs) to include as well B_1 , B_2 decay factors, ϵ , ...

→ Manual Hyperparameter tuning / Rule of Thumb Values or better → Random Search with a Coarse to Fine Scheme with appropriate ranges for each Parameter.

[CNN Sec. 11.6]

8) Training Set Validation Set Test Set

- to train the model
- To tune hyperparameters
- As a final assessment to model perf. (before deployment)
- Estimate prediction error for model selection
- Good Practice
60% - 20% - 20% split
unless dataset is gigantic then (~10%)
98% - 1% - 1%

9)

	Overfitting	Underfitting
--	-------------	--------------

Training data Performance	High	Low
---------------------------	------	-----

Testing data Performance	Low	Low
--------------------------	-----	-----

- * High Variance
(modelling training noise)
- * High Bias
(model with erroneous assumptions)

Sources of error
are bias & variance

Ex.

- e.g. Fitting an nth degree Polynomial on n+1 Points (dataset)

~~when the feature~~

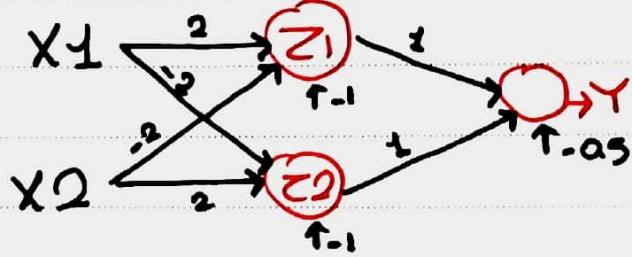
$\sum_{i=1}^n$ (instead of $\sum_{i=1}^{n-1}$)

- Fitting a linear model on highly non linear data



$$\begin{aligned} U(2x_1 - 2x_2 - 1) \\ \rightarrow \\ U(z_1 + z_2 - \frac{1}{2}) \end{aligned}$$

X1	X2	Z1	Z2	Y
0	0	0	0	0
0	1	0	1	1
1	0	1	0	1
1	1	0	0	0



- Binary Input & Output

→ The function is clearly XOR

- * If you look closely

→ Y performs OR
→ Z1 performs $X_1 \bar{X}_2$
→ Z2 performs $\bar{X}_1 X_2$

14)

$$Y = X_1 \cdot X_2 + \overline{X_1} \cdot \overline{X_2}$$

- Need Z_1, Z_2 that implement the two terms

- Observe the following

→ Suppose we wanted to implement
 $Z = X_1 \overline{X_2} \cdot X_3 \cdot X_4 \cdot \overline{X_5}$

then Clearly

$$X_1 - X_2 + X_3 + X_4 - X_5 + C$$

is equal to zero only when $C=3$ ($Z=1$)
 and $X_1 = X_3 = X_4 = 1, X_2 = X_5 = 0$. For
 any other assignment, it's less than 0. ($Z=0$)

no. of non-zeroes



Hence,

$$Z_1 = X_1 \cdot X_2 \leftrightarrow \begin{matrix} x_1 \xrightarrow{\pm 1} \\ x_2 \end{matrix} \begin{matrix} \text{circle} \\ \uparrow 1 \\ \uparrow -2 \end{matrix}$$

$x_1 + x_2 > 2$ iff &
only if $x_1 = x_2 = 1$

$$Z_2 = \overline{X_1} \cdot \overline{X_2} \leftrightarrow \begin{matrix} x_1 \xrightarrow{-1} \\ x_2 \end{matrix} \begin{matrix} \text{circle} \\ \uparrow -1 \\ \uparrow 0 \end{matrix}$$

$-x_1 - x_2 > 0$ iff &
only if $x_1 = x_2 = 0$

- * This would always work (universal and) but you could've rather thought about

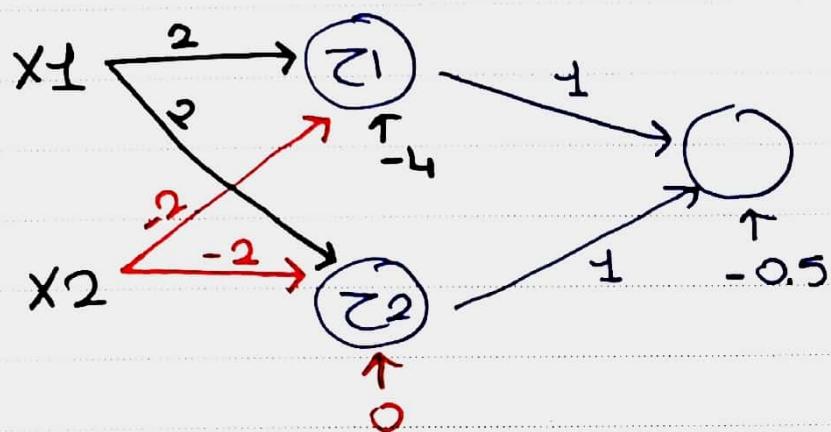
$$Z_1 = X_1 \cdot X_2 \leftrightarrow aX_1 + bX_2 + c > 0$$

a, b, c
 make this
 only odd
 For $X_1 = X_2 = 1$

- Notice that $U(x) = U(ax)$, hence always multiply all weights into a neuron by $+c$ int.

(we will
 $\times 2$)

XNOR:



12.

Classic Gradient Descent

$$w_{t+1} = w_t - \gamma \frac{\partial J}{\partial w_t}$$

stochastic,
minibatch or full
depending on how
many examples we
sum over.

Gradient Descent with momentum

$$w_{t+1} = w_t - \gamma \cdot m_t$$

where

$$m_t = \beta m_{t-1} + (1-\beta) \frac{\partial J}{\partial w_t}$$

RMS Prop

$$w_{t+1} = w_t - \frac{\gamma}{\sqrt{v_t + \epsilon}} \cdot \frac{\partial J}{\partial w_t}$$

$$v_t = \beta v_{t-1} + (1-\beta) \left(\frac{\partial J}{\partial w_t} \right)^2$$

ADAM

$$\omega_{t+1} = \omega_t - \gamma \frac{m_t}{\sqrt{v_t} + \epsilon}$$

where

$$m_t = B_1 m_{t-1} + (1-B_1) \frac{\partial J}{\partial \omega_t}$$

$$v_t = B_2 v_{t-1} + (1-B_2) \left(\frac{\partial J}{\partial \omega_t} \right)^2$$

13.

→ Gradient descent with momentum performs the update based on an exponential weighted moving average of the current & previous gradients rather than just the current gradient.

- Smoothes out oscillations as an average of gradients with emphasis on recent ones is a good estimate of the intended direction.



14. [NN Lec. 11, 9:11] (The two types of regularization)
(besides actually making the network smaller, training with more data, less features,...)

15)

1- Learning Rate → Momentum Param. → Hidden layer Nodes → Mini-batch Size

→ # Hidden layers → Learning Rate Decay → ADAM Parameters.

2- Softmax → Outputs a Probability distribution
→ only Multi-class classification
(Output layer)

3- # Hidden layers ↑ → . May be Subject to Overfitting
. Takes more time / Slower Convergence
. Network learns more complex Relations & Features*

4- ReLU recommended for Images

5- Siamese Networks → Match & Compare two Inputs

6- RNN → Has skip connections from Output to Input layer

7- Autoencoder → generate Features from Input Images and represent them in compressed low dim. Space

8- Subject to overfitting → More training data

↓ → L2 Regularization

Dropout
Regularization

9- Time Series = Trend + Seasonality + Random Noise

10- 1 Hidden layer will be enough for any boolean expression (can be written as sum of products)
OR
AND
(Input) (Hidden)