

1. Loss
2. Loss
3. ↓
4. ↓

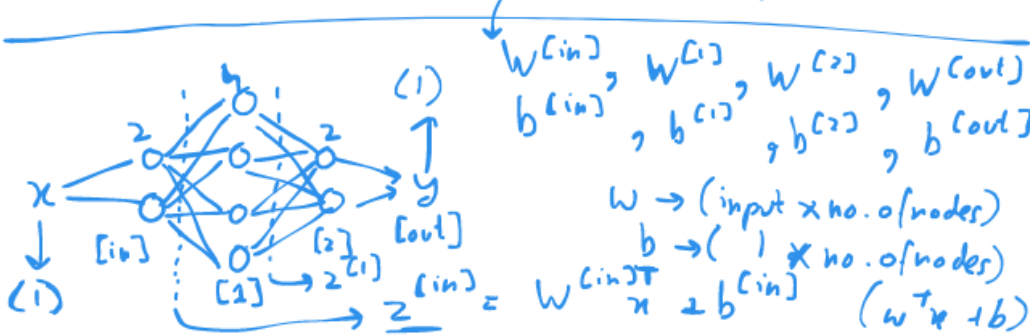
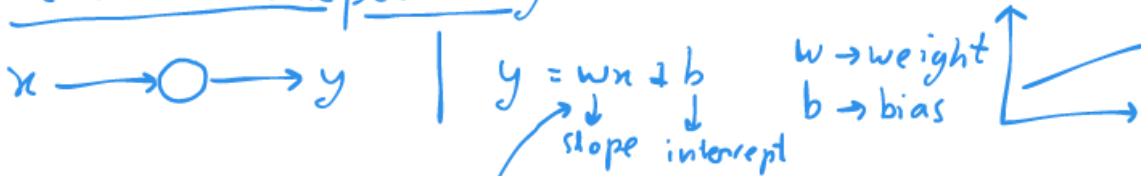
GD:

$$w' = w - \alpha \frac{\partial L}{\partial w}$$

$$\text{Loss} = (\hat{y} - y)^2$$

$\hat{y} \rightarrow \text{true}$
 $y \rightarrow \text{pred}$

Lecture 6 - Deep learning 2



$$z^{(1)} = \underbrace{W^{(1)T}}_{(2,4)} z^{(in)} + \underbrace{b^{(1)}}_{(4)} \Rightarrow (4) \rightarrow z^{(2)} \rightarrow z^{(out)} \rightarrow y$$

$$z^{(1)} = [w_1, w_2] \cdot x + [b_1, b_2]$$

$$z^{(1)} = \boxed{W^{(1)T} (W^{(in)T} x + b)} + b^{(1)}$$

always linear

Activations \rightarrow Non-linearity

\downarrow
 Approximate Complex Data
 (for eg. Images, Audio, Robot Joint Data)

$$x = \begin{bmatrix} [-1] \\ [-0.400] \\ [-0.999] \\ \vdots \end{bmatrix} \quad (b \times 1)$$

1. Loss
2. Loss
3. Loss
4. Loss

$$\text{Loss} = (\hat{y} - y)^2$$

$\hat{y} \rightarrow \text{true}$
 $y \rightarrow \text{pred}$

GD:

$$w' = w - \alpha \frac{\partial L}{\partial w}$$

\rightarrow Slow down

Lecture 6 - Deep learning 2



$$z = wx + b$$

$w \rightarrow \text{slope}$ $b \rightarrow \text{intercept}$

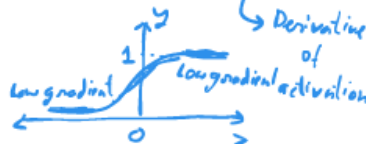
$w \rightarrow \text{weight}$
 $b \rightarrow \text{bias}$

$$\frac{\partial L}{\partial w} \rightarrow \frac{\partial L}{\partial y} \left(\frac{\partial y}{\partial z} \right) \frac{\partial z}{\partial w}$$

Derivative of activation

✓ ① Sigmoid

$$\sigma(z) = \frac{1}{1 + e^{-z}}$$



✓ ② Tanh

$$\tanh(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}}$$



Tanh better sigmoid

Sigmoid \leftarrow

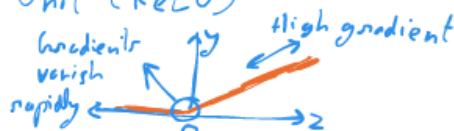
\rightarrow Better for binary

Tanh

\rightarrow Better for -ve to +ve range

③ Rectified Linear Unit (ReLU)

$$y = \max(0, z)$$



Tips:

$\sigma(z) \rightarrow$ Use only output layer

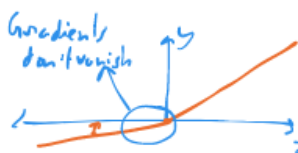
$\tanh(z)$ / ReLU \rightarrow Use for hidden layer



④ Leaky ReLU:

$$y = \max(0.01z, z)$$

\downarrow
change



⑤ ... much more \rightarrow Explore

Lecture 6 - Deep learning 2

$$x \rightarrow \bigcirc \rightarrow y \quad \bigg| \quad z = wx + b \quad \begin{matrix} w \rightarrow \text{weight} \\ b \rightarrow \text{bias} \end{matrix}$$

slope intercept



\rightarrow Not enough parameters

\rightarrow Not long enough training

\rightarrow Get more unique data

\rightarrow Regularization

\hookrightarrow Fixes the variance

① L2 Regularization

$$\sum_{i=1}^m \frac{1}{m} L(\hat{y}_i, y_i) \Rightarrow \text{Loss} \quad m = \text{no. of samples}$$

$$\sum_{i=1}^m \frac{1}{m} L(\hat{y}_i, y_i) + \frac{\lambda}{2m} \|w\|^2$$

$\lambda \Rightarrow$ scalar

$\frac{\lambda}{2m} \|w\|^2$ \Rightarrow L2 norm $\sqrt{\sum_i w_i^2}$

\rightarrow Detects Menor

\rightarrow Detects Femoris

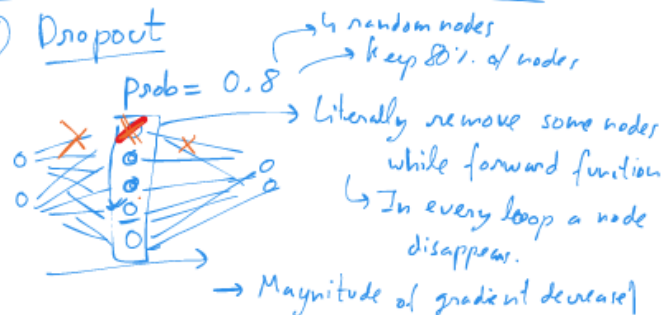
\rightarrow learn abt derivation $w_3 \gg w_1$ \rightarrow Model overfits to detecting Femoris

Weighted squared sum of weights

\rightarrow Added in loss

\hookrightarrow i.e. Try to reduce very high values of weights

② Dropout



- Not enough parameters
- Not long enough training
- Get more unique data
- Regularization

→ Fixes the variance

① L2 Regularization

$\sum_{i=1}^m \frac{1}{m} L(\hat{y}_i, y_i) \Rightarrow \text{loss}$ $m = \text{no. of samples}$

$\sum_{i=1}^m \frac{1}{m} L(\hat{y}_i, y_i) + \frac{\lambda}{2m} \|w\|^2$

→ Detects Menoroder
→ Detects Ferranis
→ Learn abt derivation $w_3 \gg w_1$ → Model overfits to detecting Ferranis

$\lambda \Rightarrow \text{Scalar}$

Regularization term

$[w_1, w_2, w_3, w_4]$
↓ L2 norm
 $\sqrt{\sum_i (w_i)^2}$

Weighted squared sum of weights
→ Added in loss
→ ∴ Try to reduce very high values of weights

② Dropout

