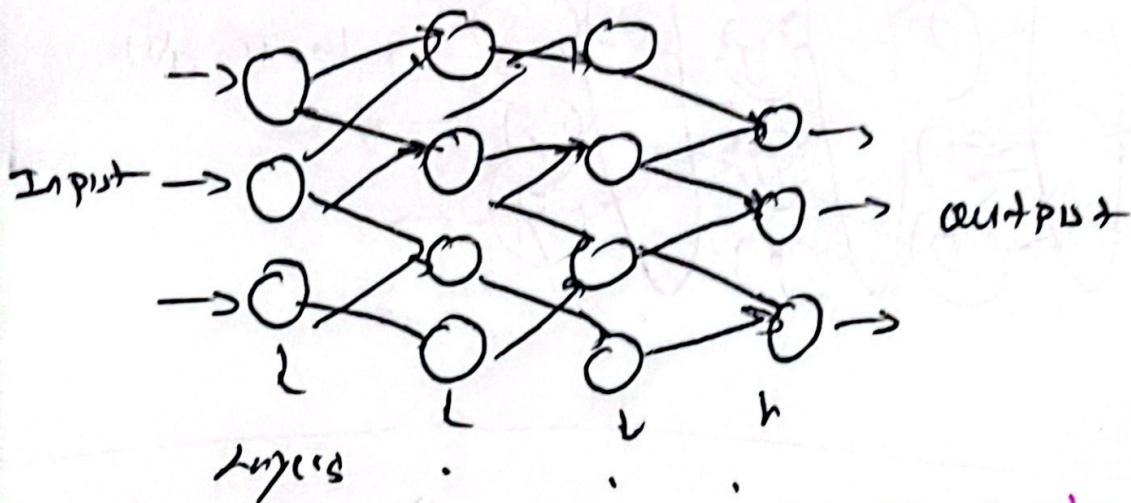


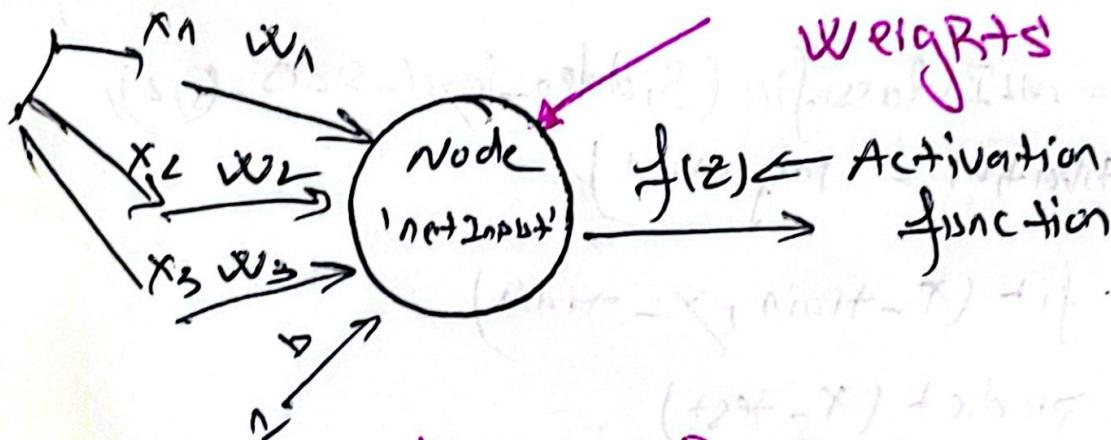
① Neural Net Structure By Limbo



Data

Combined via

Weights



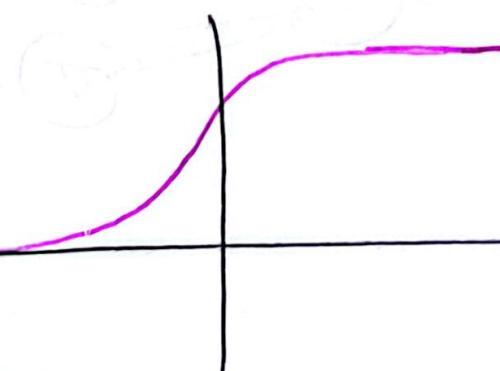
↪ Relation to Logistic Regression

$$f(z) = \frac{1}{1 + e^{-z}} \quad z = b + \sum_{i=1}^m x_i w_i + b$$

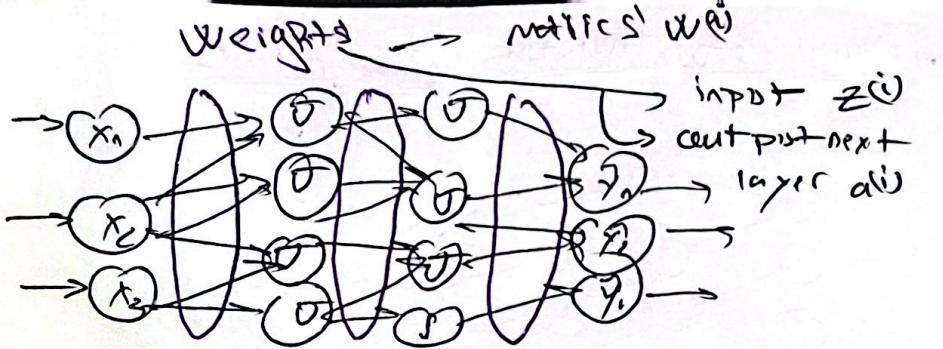
→ Sigmoid function

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

$$\sigma'(z) = \sigma(z)(1 - \sigma(z))$$



①



$$z^{(1)} = \alpha^{(1)} W^{(1)} \rightarrow c$$

$$z_i^{(2)} = \alpha^{(2)} W^{(2)} \rightarrow$$

$$z^* = \alpha^{(n-1)} W^{(n)}$$

→ Input + training input
→ Compute \rightarrow
→ Adjust and recompute
→ Back propagation
→ Output $\rightarrow f_1$
 \circ
 \circ

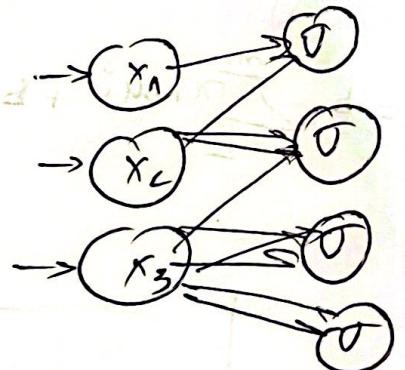
code

```
from sklearn.neural_network import MLPClassifier
```

```
m1P = MLPClassifier(hidden_layer_sizes=(5, 2),  
activation='logistic')
```

```
m1P.fit(X_train, Y_train)
```

```
m1P.predict(X_test)
```



$W^{(1)}$ 3x4 matrix

$z^{(1)}$ uvector

α^2 uvectorial

$$\alpha = \alpha^0$$

$$z^* = \alpha^{(n-1)} W^{(n)}$$

$$\alpha^* = \sigma(z^*)$$

$$C \rightarrow \frac{\delta E}{\delta w_k}$$

④

$$z^{(1)} = x w^{(1)} \rightarrow d^{(1)} = J(z^{(1)})$$

$$z^{(2)}_i = a^2 w^{(2)} \rightarrow a^2 = J(z^{(2)})$$

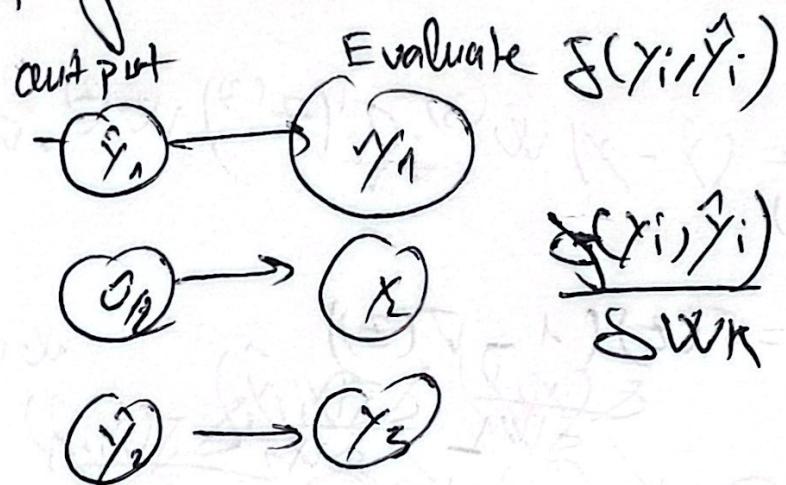
$$z^n = d^{n-1} w^{(n)} \rightarrow \hat{x} = \text{softmax}(z^n)$$

→ Input training inputs → True output

→ Compute $\underset{\text{loss function}}{\longrightarrow}$

→ Adjust and repeat

→ Back propagation → Run to make adjustment



$$\therefore \frac{\delta J}{\delta w_k}$$

(2)

Back propagation : Feed forward network

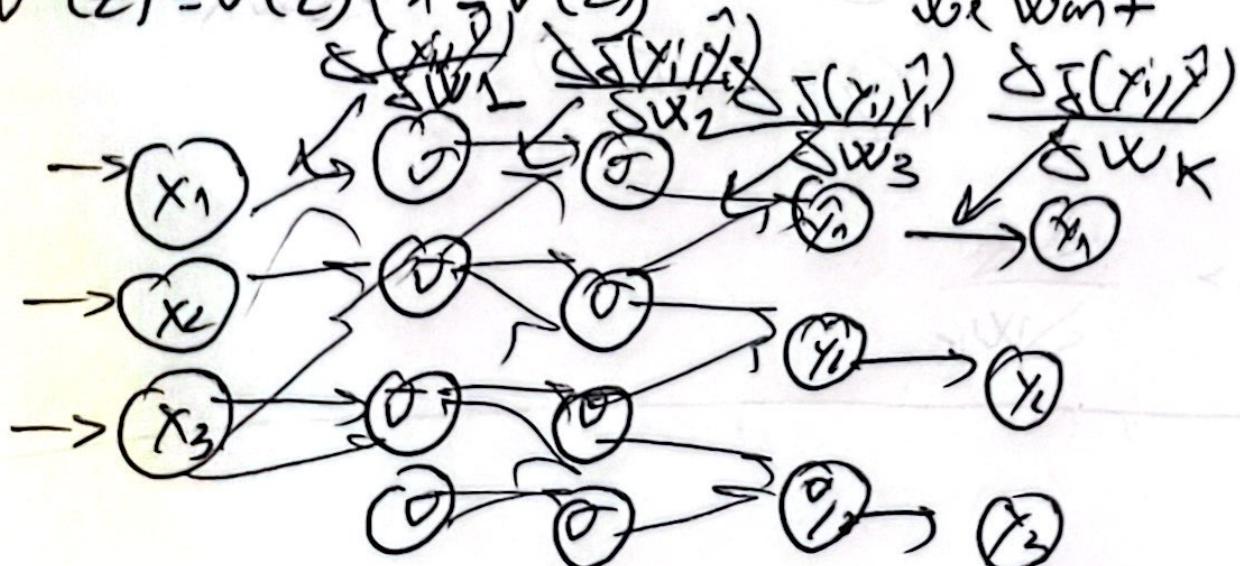
$$\frac{\partial E}{\partial w^{(n)}} = (\hat{y} - y) \alpha^{(n)}$$

$$\frac{\partial E}{\partial w^{(3)}} = (\hat{y} - y) \alpha^{(3)}$$

$$\frac{\partial E}{\partial w^{(2)}} = (\hat{y} - y) \cdot w^{(3)} \cdot \sigma'(z^{(2)}) \alpha^{(1)}$$

$$\frac{\partial E}{\partial w^{(1)}} = (\hat{y} - y) w^{(3)} \cdot \sigma'(z^{(2)}) w^{(2)} \cdot \sigma'(z^{(1)}) \alpha^{(0)}$$

$$\sigma'(z) = \sigma(z)(1 - \sigma(z))$$



(3)

→ mathematical understanding

$$\hat{y} = f(u_0x_0 + \dots + u_nx_n + b)$$

→ loss function

$$L = \frac{1}{2} (y - \hat{y})^2$$

$$\frac{\delta L}{\delta w_i} = \frac{\delta L}{\delta y} \frac{\delta y}{\delta w_i}$$

$$= -(y - \hat{y})$$

$$w = w - \eta \frac{\delta L}{\delta w}$$

→ vanish gradient

$$z^{(l)} = w^{(l)} a^{(l-1)} + b^{(l)}$$

⑥

$$a^{(l)} = f(z^{(l)})$$

$$\frac{\delta L}{\delta w_i} = \prod_j f'(z^{(j)}) w^{(j)}$$

$$f'(z) = f(z)(1-f(z))$$

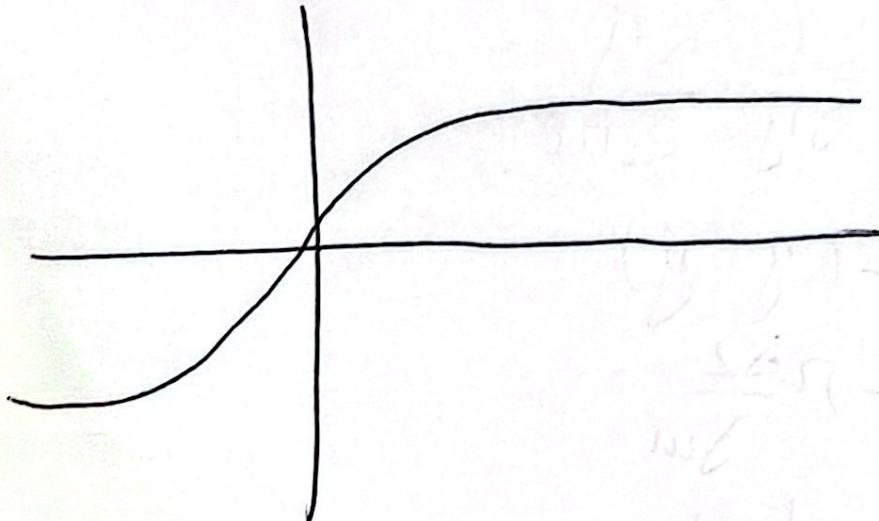
$$0 < f'(z) < 0.25$$

Grad $\rightarrow 0 \rightarrow$ how to fix
activate ReLU $f(z) = \max(0, z)$

Raz deivate 1 when $z > 1$

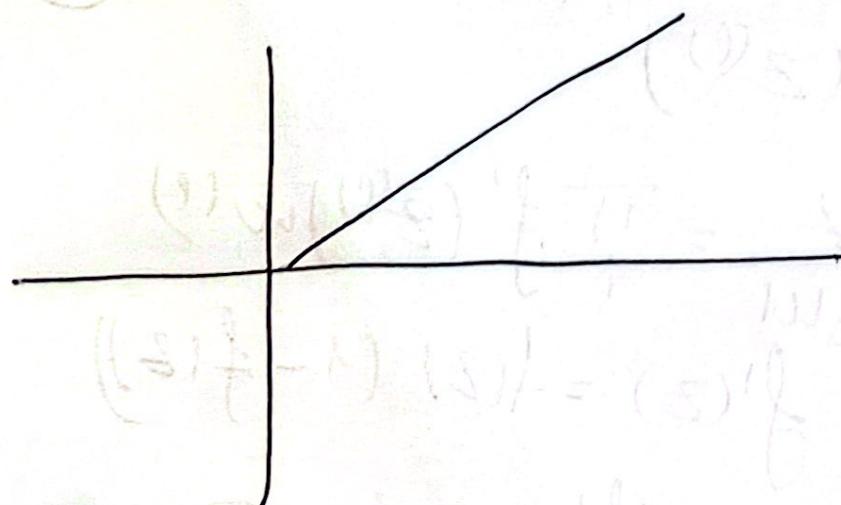
→ Activation function

$$\tanh(z) = \frac{\sinh(z)}{\cosh(z)} = \frac{e^{2x} - 1}{e^{2x} + 1}$$



(ReLU) →

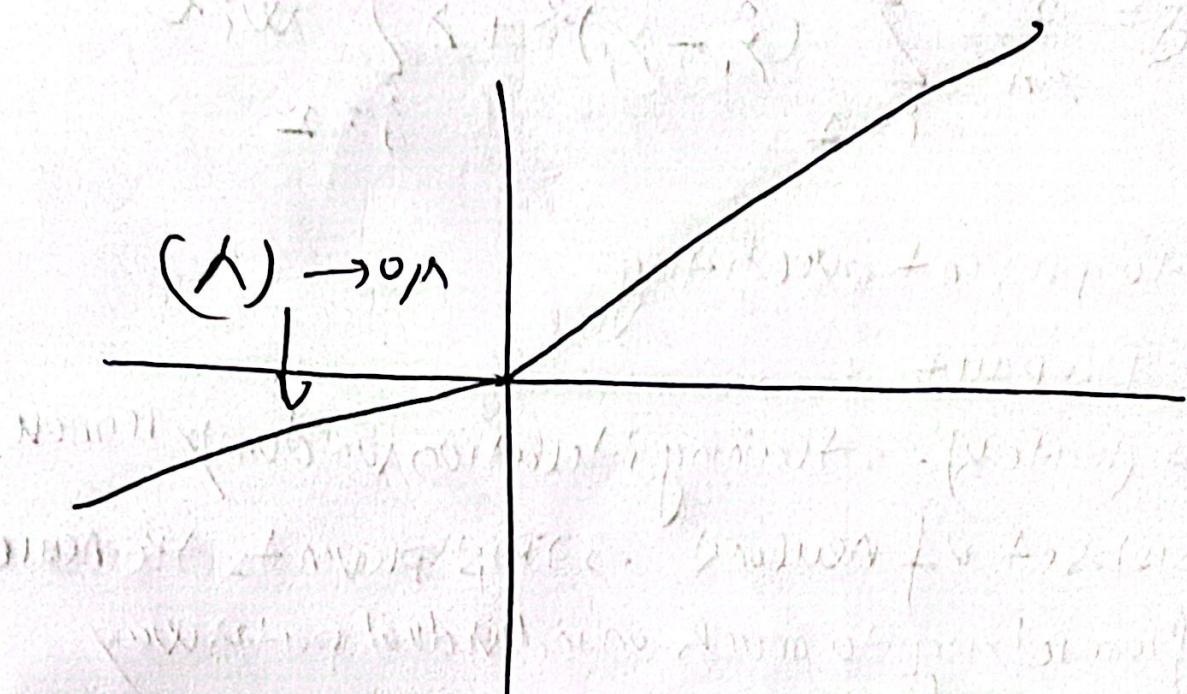
$$\text{ReLU}(z) \begin{cases} 0 & z < 0 \\ z, & z \geq 0 \end{cases} = \max(0, z)$$



②

(Leaky) ReLU

$$\text{LeReLU}(z) = \begin{cases} \alpha z & z < 0 \\ z & z \geq 0 \end{cases} \quad \max(\alpha z, z) \text{ for } \alpha < 1$$



Optimizers

$$\min_w \text{Loss}(w)$$

$$w = w - \alpha \frac{\text{Loss}}{\partial w} \rightarrow \text{learning Rate}$$

$$\hookrightarrow \text{GD} \quad w = w - \alpha \nabla \text{Loss}(x_i)$$

Converge a lot

(A)

(B)

Regularization

Regularization → Reduce error

Cost-function

$$J = \frac{1}{2n} \sum_{i=1}^n (\hat{y}_i - y_i)^2 + \lambda \sum_{j=1}^m w_j^2$$

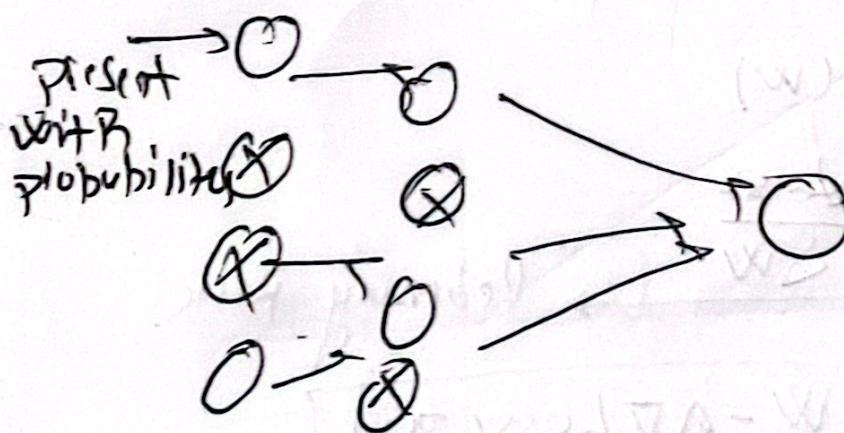
→ to prevent overfitting

→ Dropout

↪ (batch) → training iteration (randomly remove

subset of neurons) → This prevents the neuron from relying too much on individual pathway

make it more "robust"



(g)

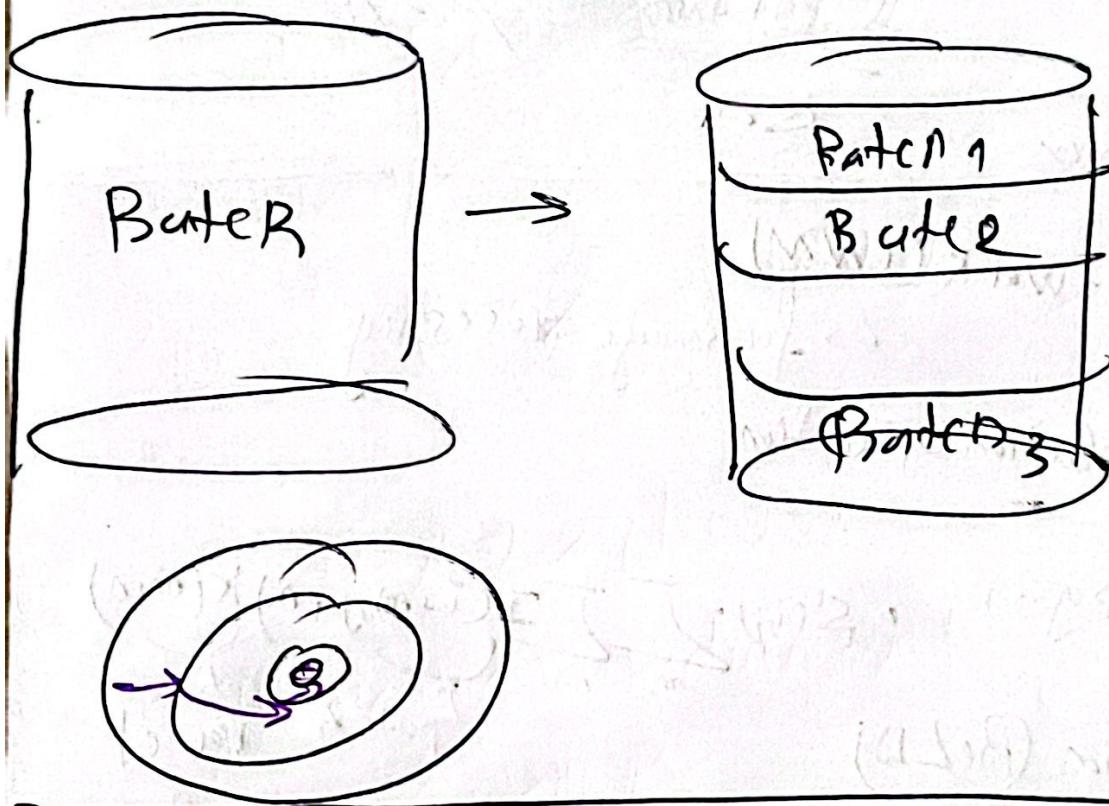
Data Shuffling

↳ Randomizing + Reorder

↳ prevent overfitting

↳ Improve (SGD)

↳ Reduce Bias



→ for $(0, 1)$ classification problems → Sigmoid activation \equiv Logistic Regression

1 + Take a vector with length \rightarrow number of categories

2 + Category \equiv (1)

Creating $\begin{bmatrix} 1 \\ 0 \end{bmatrix}$

saving $\begin{bmatrix} 0 \\ 1 \end{bmatrix}$

Mortgage $\begin{bmatrix} 0 \\ 1 \end{bmatrix}$

10

Wiederholung

$$\frac{\partial \hat{y}_i}{\partial x_i} = \frac{\partial}{\partial x_i} (w_i x_i + b) = w_i$$

Gradient

▼ Visualized

⑥ Fully connected layer

③ Pooling \rightarrow Reducing the Dimension

$$f(x) = \max(0, x)$$

Activation (ReLU)

$$(Imag(-)) \rightarrow f(i, j) = \max(0, i + m_j + b)$$

Pooling (Max)

④ Convolution operation

\hookrightarrow for image processing

Neural Network (CNN)

$$CE = -\frac{\log p_{\text{true}}}{\text{Softmax}}$$

$$CE = -\sum_i p_i \log(p_i)$$

for loss function we use "CE"

$$\frac{\partial CE}{\partial z_i} = \frac{\partial}{\partial z_i} (-\sum_i p_i \log(p_i))$$

"Fine-Tuning" \Rightarrow pre-trained network



① Pretrained Model

$$f_{\theta}(x) = y$$

↳ model parameter

$$\theta^* \rightarrow \underline{\text{optimized}}$$

$$\theta^* = \arg \min \sum_{(x,y) \in D} \mathcal{L}(f_{\theta}(x), y)$$

C optimization

$$\theta_{\text{fine-tune}}^* = \arg \min \sum_{(x,y) \in D} \mathcal{L}(f_{\theta}(x), y)$$

C SGD

$$\theta \leftarrow \theta - \eta \nabla \mathcal{L}(f_{\theta}(x), y)$$



② Regularization to Reduce overfitting

$$\mathcal{L}_{\text{new}} = \mathcal{L} + \lambda \|\theta\|^2$$

③

$$f_{\theta^*}(x) \approx f_{\theta_{\text{new}}}(x)$$

$$\theta_{\text{fin}} \approx \theta + \Delta \theta$$

RNN, RL, GAN

12

(RNN) → for sequential data (time series)

text, speech) → allowing info persistence

$$h_t = \sigma(W_R h_{t-1} + W_X x_t + b_R)$$

$$y_t = \sigma(W_y h_t + b_y)$$

Application → language modeling
↳ machine translation
↳ speech rec

2) RL

$$V^\pi(s) = E_\pi \left[\sum_{k=0}^{\infty} \gamma^k r_{t+k} | s_t = s \right]$$

$$\gamma \in [0, 1]$$

$$\pi(a, s') \text{ probability for } \xrightarrow{s'} a$$

:

Pedda abzione

