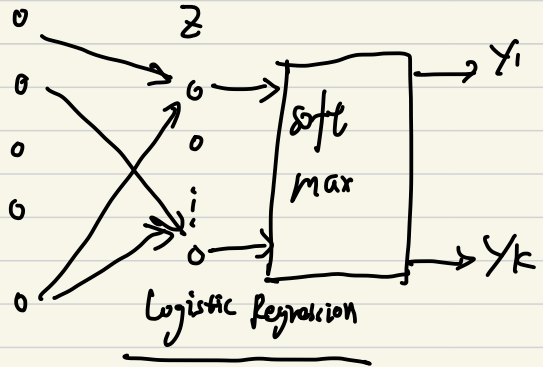
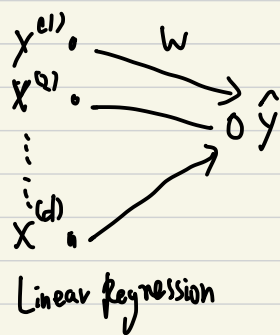


L9



- Feed Forward
- Directed acyclic graphs

Neural Network

primitive $z = \phi(\langle w, x \rangle + b)$
"Neuron"

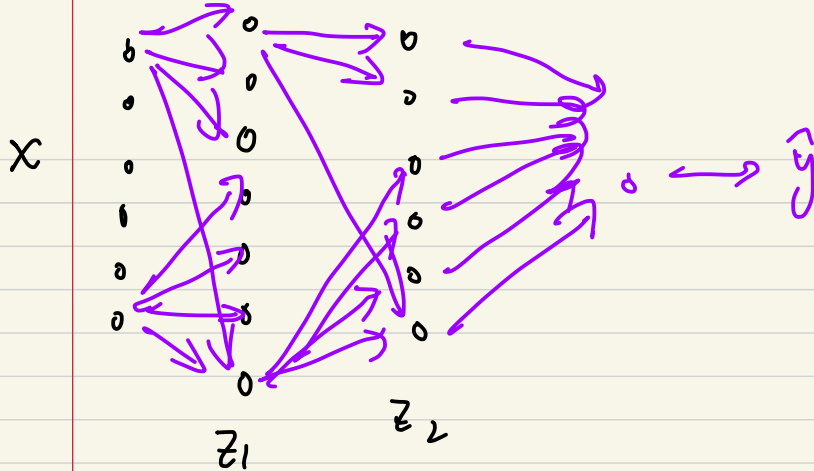


ϕ "activation function"

Examples: $\phi(z) = z \rightarrow$ linear regression

$\phi(z) = \frac{1}{1 + e^{-\langle w, x \rangle + b}} \rightarrow$ logistic regression

$\phi(z) = \text{sign}(z) \rightarrow$ perceptron



$$W^1 \quad W^2 \quad W^3$$

$$z_1 = \sigma^{(1)}(W^{(1)}x + b^{(1)})$$

$$z_2 = \sigma^{(2)}(W^{(2)}z_1 + b^{(2)})$$

$$\hat{y} = \sigma^{(3)}(W^{(3)}z_2 + b^{(3)})$$

" 3-Layer network / 2-hidden layer network

Why chosen? — historical / biological

Universal approximation Theorem (Cybenko, 1988)

≈ any continuous function can be approximated by
 a 1-hidden layer network of finite size (for most
 functions, activation).

Drawback

- Existence Theorem only
- Practical Issue

Representation

2 What network size & shape

Training

2 How to choose the network weights

Generalization

2 Does the learned generalize to unseen input?

Neural network architecture

Design criteria : 1) Weights

2) Activation function

Activation Functions:

this is linear $\Rightarrow \bullet \sigma(z) = z$

\Downarrow [Not too useful for multi-layer networks]

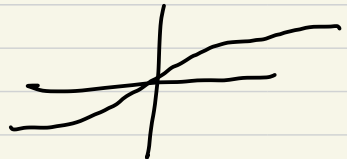
$$y = W^{(3)} W^{(2)} W^{(1)} X$$

$$= \bar{W} X$$

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

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

(tanh)

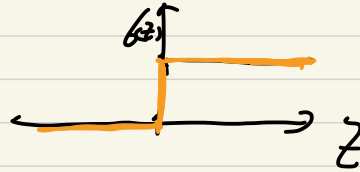


- $\sigma(z) = \text{ReLU}(z)$
 $= \max(z, 0)$



- $\sigma(z) = \text{sgn}(z)$

- $\sigma(z) = \text{HT}(z)$



Weights:

- Dense layers \rightarrow weights are arbitrary
- Convolution layers \rightarrow weights implement convolution
- Pooling layers \rightarrow Downsampling size of output
- Batch normalization \rightarrow Rescaling
- Recurrent layers \rightarrow Feed back
- Attention layers \rightarrow NLP (etc.)

Q: How do I mix & match?

A: No correct answer

Thumb Rule: Just use a good existing architecture

Training

Given architecture, how do I train a network?

★ Define loss function (+ Optional Regularization)

★ Use some variant of gradient descent (GD/SGD/...)

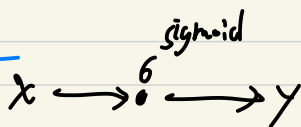
$$\{(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)\}$$

$$L(w) = \frac{1}{n} \sum_{i=1}^n \ell(y_i, f_w(x_i)) + \lambda R(w)$$

$$w_{t+1} \leftarrow w_t - \alpha^t \nabla L(w)$$

Hard to compute

Toy example



$$z = wx + b$$

$$f(z) = \sigma(z)$$

$$L(w, b) = \frac{1}{2} (y - f(x))^2 + \lambda w^2$$

$$= \frac{1}{2} (y - \sigma(wx + b))^2 + \lambda w^2$$

$$\nabla L(w, b) = \begin{pmatrix} \frac{\partial L}{\partial w} \\ \frac{\partial L}{\partial b} \end{pmatrix}$$

$$\frac{\partial L}{\partial w} = (\sigma(wx + b) - y) \cdot \sigma'(wx + b) \cdot x + 2\lambda w$$

$$\frac{\partial \mathcal{L}}{\partial b} = (\delta(wx+b) - y) \delta'(wx+b)$$

it's 4' efficiency

Back propagation / back-prop.

- exploit the structure of the network

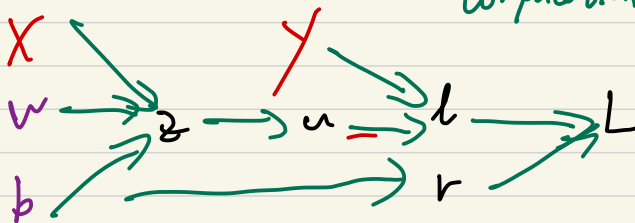
$$z = wx + b$$

$$u = \delta(z)$$

$$l = \frac{1}{2} (u - y)^2$$

$$r = u^2$$

$$L = l + \lambda r$$



Computation graph for a single neuron.

FORWARD (graph) (value)

For node $i = 1, 2, \dots, N$

compute V_i as a function of parents of i

Backward (graph) : gradient (前向传播)

For node $i = N, \dots, 1$:

$$\text{compute } \frac{\partial \mathcal{L}}{\partial V_i} = \sum_{j \in \text{children}(i)} \frac{\partial \mathcal{L}}{\partial V_j} \cdot \frac{\partial V_j}{\partial V_i}$$

$i = 1, 2, \dots, N$
 $N \leftarrow$ number of nodes

$$\frac{\partial L}{\partial v} := \partial_v L \text{ (表的形式)}$$

$$\partial_v L := 1$$

$$\partial_v L := 1$$

$$\partial_v L = \lambda$$

chain rule

$$\frac{\partial L}{\partial u} := \partial_u L \cdot \partial_u u = 1 \cdot (u - y)$$

$$\partial_x L = \partial_u L \cdot \partial_z u$$

$$= \partial_u L \cdot \phi'(z)$$

$$\partial_w L = \partial_x L \partial_w x + \partial_f L \partial_w y$$

$$\partial_b L = \partial_f L \cdot \partial_b z$$

$$= \partial_z L$$

没有重复计算

不像会计算那率

一步步列出来重

复代值

• Computationally efficient

• Modular 一个个模块, 很标准, 可以套用

• Scalable 适用于更大规模