

Types of NNs

Wen Yu

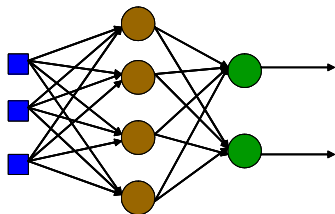
<https://www.ctrl.cinvestav.mx/~yuw/>
yuw@ctrl.cinvestav.mx

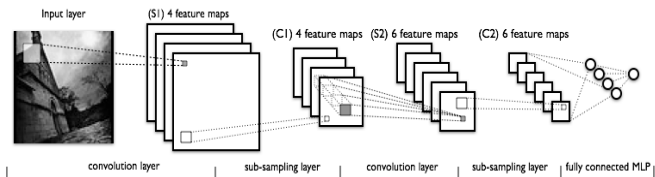
Types of NNs

- Time domains
 - ① Discrete-time
 - ② Continuous-time
- Structure
 - ① Feedforward
 - ② Recurrent

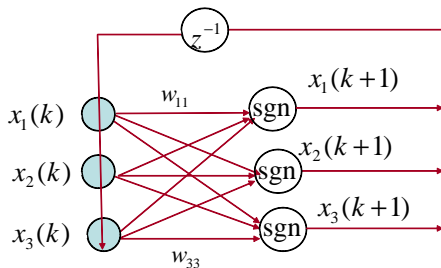
So

- ① Discrete-time Feedforward: MLP
- ② Discrete -time Recurrent: RNN
- ③ Continuous-time Feedforward: MLP
- ④ Continuous-time Recurrent: *Differential Neural Networks*





Parallel form



- Boltzmann machines
- Long-short term memory (LSTM)

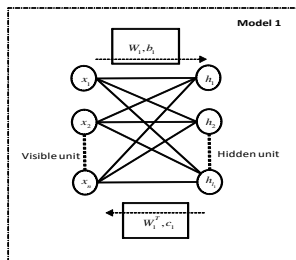
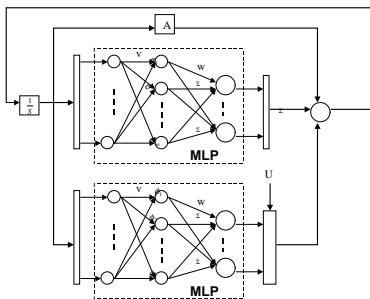


Figure: Markov sampling in a restricted Boltzmann machine

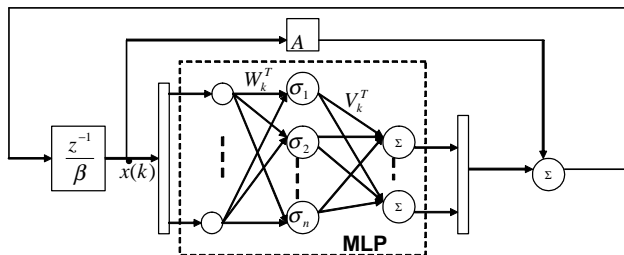
Alexander S Poznyak, Edgar N Sanchez & Wen Yu , *Differential Neural Networks for Robust Nonlinear Control*, World Scientific Publishing Co., 2001

Dynamic NNDifferential NN



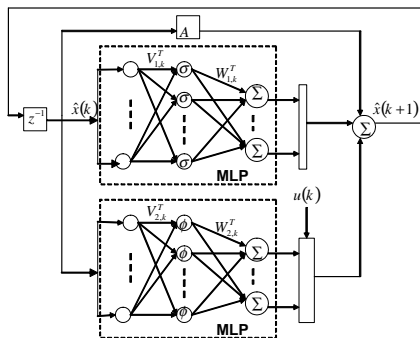
$$\dot{x}_t = Ax_t + W_{1,t}\sigma(V_{1,t}x_t) + W_{2,t}\phi(V_{2,t}x_t)u_t$$

Discrete-time recurrent neural networks



$$\beta \hat{x}(k+1) = A \hat{x}(k) + V_k \sigma [W_k x(k)]$$

Discrete-time recurrent neural networks



$$\hat{x}(k+1) = A\hat{x}(k) + V_{1,k}\sigma[W_{1,k}x(k)] + V_{2,k}\sigma[W_{2,k}x(k)]u_k$$

Universal approximation of *Differential Neural Networks*

We use the Differential Neural Network

$$\frac{d}{dt}\hat{x} = A\hat{x} + W\Phi[V\hat{x}]$$

to model a nonlinear system, $A \in R^{n \times n}$, $W \in R^{n \times m}$, $V \in R^{m \times n}$, $x \in R^n$

$$\dot{x} = f(x)$$

Theorem

Suppose the dynamic system and the differential Neural Network are initially at the same state $x_0 = \hat{x}_0$. For any $\varepsilon > 0$ and any finite $T > 0$, there exists an integer m such that the state of differential Neural Network \hat{x}_t with the weight $[W^, V^*]$ satisfies*

$$\sup_{0 \leq t \leq T} |x_t - \hat{x}_t| < \varepsilon \quad (1)$$

Universal approximation of *Differential Neural Networks*

The nonlinear system, can be expressed as

$$\dot{x} = Ax_t + f(x) - Ax_t = Ax_t + G(x)$$

where $G(x) = f(x) - Ax$

From Stone-Weierstrass theorem if the hidden layer V is large enough, the neural network with one hidden layer

$$\hat{y} = W\Phi[V\hat{x}]$$

can approximate any

$$y = G(x)$$

to any degree of accuracy,

$$\sup_x |G(x) - W\Phi[Vx]| \leq \delta \quad (2)$$

where δ is any small positive constant.

Wen Yu, Nonlinear system identification using discrete-time recurrent neural networks with stable learning algorithms, Information Sciences, Vol.158, No.1, 131-147, 2004.

Let us define identification error as

$$\Delta_t = x_t - \hat{x}_t \quad (3)$$

There exists an integer m and the weight $[W^*, V^*]$ such that

$$\dot{\hat{x}}_t = A\hat{x}_t + W^* \Phi [V^* \hat{x}]$$

because

$$\dot{x} = Ax_t + G(x)$$

so

$$\dot{\Delta}_t = A\Delta_t + G(x) - W^* \Phi [V^* \hat{x}] \quad (4)$$

Since $x_0 = \hat{x}_0$, $\Delta_0 = 0$. The solution of the differential equation

$$\dot{\Delta}_t = A\Delta_t + G_1$$

where $G_1 = G(x) - W^*\Phi[V^*\hat{x}]$

$$\begin{aligned} \Delta_t &= \int_0^t e^{A(t-\tau)} G_1 d\tau \\ &= \int_0^t e^{A(t-\tau)} [W^*\Phi[V^*x] - W^*\Phi[V^*\hat{x}]] d\tau \\ &\quad + \int_0^t e^{A(t-\tau)} [G(x) - W^*\Phi[V^*x]] d\tau \end{aligned} \quad (5)$$

Since A is a stability matrix, there exists a positive constant α such that

$$\|e^{At}\| \leq e^{-\alpha t}, \quad \alpha > 0$$

So

$$\begin{aligned} |\Delta_t| &\leq \int_0^t \|e^{A(t-\tau)}\| \|W^*\| \{|\Phi[V^*x] - \Phi[V^*\hat{x}]\}| d\tau \\ &\quad + \int_0^t \|e^{A(t-\tau)}\| \{|G(x_t, u_t) - W^*\Phi[V^*x]|\} d\tau \end{aligned}$$

Using Lipschitz condition

$$|\Phi[V^*x] - \Phi[V^*\hat{x}]| \leq l|x_\tau - \hat{x}_\tau|$$

and Stone-Weierstrass theorem

$$\sup_{(\hat{x}, u)} |G(x) - W^*\Phi[V^*x]| \leq \delta \tag{6}$$

So

$$\begin{aligned}
 |\Delta_t| &\leq \int_0^t e^{-\alpha(t-\tau)} L |x_\tau - \hat{x}_\tau| d\tau + \int_0^t e^{-\alpha(t-\tau)} \delta d\tau \\
 &= \int_0^t e^{-\alpha(t-\tau)} L |\Delta_t| d\tau + \frac{\delta}{\alpha} (1 - e^{-\alpha t})
 \end{aligned}$$

where $L := I \|W^*\|$. Because $e^{-\alpha t} \leq 1$, for $\alpha t > 0$

$$|\Delta_t| \leq \frac{\delta}{\alpha} + \int_0^t e^{-\alpha(t-\tau)} L |\Delta_t| d\tau$$

Proof

For

$$|\Delta_t| \leq \frac{\delta}{\alpha} + \int_0^t L e^{-\alpha(t-\tau)} |\Delta_t| d\tau$$

We use Gronwall-Bellman inequality, if

$$u \leq u_0 + \int f(\tau) u(\tau) d\tau$$

then

$$u \leq u_0 e^{\int f(\tau) d\tau}$$

But

$$u_0 = L \frac{\delta}{\alpha} \int_0^t (1 - e^{-\alpha t}) e^{-\alpha(t-\tau)}$$

$$\begin{aligned} |\Delta_t| &\leq \frac{\delta}{\alpha} + L \frac{\delta}{\alpha} \int_0^t (1 - e^{-\alpha t}) e^{-\alpha(t-\tau)} e^{\int_{\tau}^t L e^{-\alpha(t-s)} ds} d\tau \\ &= \frac{\delta}{\alpha} + L \frac{\delta}{\alpha} \int_0^t (1 - e^{-\alpha t}) e^{-\alpha(t-\tau)} e^{\frac{L}{\alpha}(1 - e^{-\alpha(t-\tau)})} d\tau \end{aligned}$$

Because $e^{-\alpha t} \leq 1$, for $\alpha t > 0$.

Using the inequality of

$$1 - e^{-x} \leq x, \quad x \geq 0$$

we have

$$\begin{aligned}
 |\Delta_t| &\leq \frac{\delta}{\alpha} + L \frac{\delta}{\alpha} \int_0^t (1 - e^{-\alpha\tau}) e^{-\alpha(t-\tau)} e^{L(t-\tau)} d\tau \\
 &= \frac{\delta}{\alpha} + L \frac{\delta}{\alpha} \int_0^t (1 - e^{-\alpha\tau}) (1 - e^{-\alpha\tau}) e^{(L-\alpha)(t-\tau)} d\tau \\
 &= \frac{\delta}{\alpha} + \frac{(\alpha-L)e^{-\alpha t} - \alpha e^{-(\alpha-L)t} + L}{\alpha(\alpha-L)} \delta \\
 &= \frac{\delta}{\alpha} + \left[\frac{1}{\alpha} e^{-\alpha t} - \frac{e^{-(\alpha-L)t}}{(\alpha-L)} + \frac{L}{\alpha(\alpha-L)} \right] \delta \\
 &\leq \frac{\delta}{\alpha} + \left[\frac{1}{\alpha} - \frac{e^{(\alpha-L)t}}{(\alpha-L)} + \frac{L}{\alpha(\alpha-L)} \right] \delta \\
 &= \frac{\delta}{\alpha} + \frac{\delta}{\alpha-L} \left[1 - e^{-(\alpha-L)t} \right] \\
 &\leq \frac{\delta}{\alpha} + \frac{\delta}{\alpha-L} = \frac{(2\alpha-L)}{\alpha(\alpha-L)} \delta
 \end{aligned}$$

If we define

$$\varepsilon := \frac{(2\alpha - L)}{\alpha(\alpha - L)}\delta$$

and select $\alpha \geq L$, so $\varepsilon > 0$, $|\Delta_t| = |x_t - \hat{x}_t| \leq \varepsilon$. From function approximate theorem (2) we know δ can be any small. Because α and L are constants, so ε can also be any small.

If the initial condition of neural networks and dynamic system are different, $x_0 \neq \hat{x}_0$, The solution of the differential equation (4) is

$$\Delta_t = \int_0^t e^{A(t-\tau)} [G(x_\tau, u_\tau) - W^* \Phi_{V^*}(\hat{x}_\tau, u_\tau)] d\tau + \Delta_0 e^{At}$$

The additional term $\Delta_0 e^{At}$ will influence the final result as

$$|\Delta_t| \leq \frac{(2\alpha - L)}{\alpha(\alpha - L)} \delta + \frac{1}{L} e^{-(\alpha-L)t} |\Delta_0|.$$

Because $(\alpha - L)t > 0$, l and Δ_0 are constants

$$\lim_{t \rightarrow \infty} \frac{1}{L} e^{-(\alpha-L)t} |\Delta_0| = 0$$

So the influence of initial condition of RMLP will decay exponentially to zero.

There are three broad types of learning:

- Supervised Learning, learning with a teacher)
 - - 1 Classification
 - 2 Control
 - 3 Function approximation
 - 4 Associative memory
- Unsupervised learning, learning without help)
 - 1 Clustering
- Reinforcement learning, learning with limited feedback

Learning in biological systems as optimization

- Learning = learning by adaptation
- The young animal learns that the green fruits are sour, while the yellowish/reddish ones are sweet. The learning happens by adapting the fruit picking behavior. The animal likes to eat many energy rich, juicy fruits that make its stomach full, and makes it feel happy.
- At the neural level the learning happens by changing of the synaptic strengths, eliminating some synapses, and building new ones

Learning principle for artificial neural networks

- Maintaining synaptic strength needs energy, it should be maintained at those places where it is needed, and it shouldn't be maintained at places where it's not needed
- ENERGY MINIMIZATION
- We need an appropriate definition of energy for artificial neural networks, and having that we can use mathematical optimization techniques to find how to change the weights of the synaptic connections between neurons.
- ENERGY = measure of task performance error

Unsupervised Learning

- ANN adapts weights to cluster input data
- Hebbian learning
 - Connection stimulus-response strengthened (Hebbian)
- Competitive learning algorithms
 - Kohonen & ART
 - Input weights adjusted to resemble stimulus

Shifting around Learning

- we can think of learning as the process of shifting around the optimal points until each training data is classified correctly.
- we need to formalize that process of “shifting around” into a systematic algorithm that can easily be implemented on a computer.

$$w(k+1) = w(k) + \Delta w(k)$$

The “shifting around” can conveniently be split up into a number of small steps.

Hebbian learning

Hebbian learning (Unsupervised learning)

$$w(k+1) = w(k) + \eta \hat{y}(k) u(k)$$

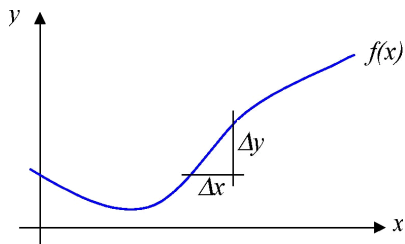
Perceptrons learning (supervised learning)

$$w(k+1) = w(k) + \eta [y(k) - \hat{y}(k)] u(k)$$

There is some similarity, but it is clear that Hebbian learning is not going to get our Perceptron to learn a set of training data. There are variations of Hebbian learning that do provide powerful learning techniques for biologically plausible networks, such as Contrastive Hebbian Learning, but we shall adopt another approach for formulating.

Delta Rule

We want to change the value of x to minimise $f(x)$



Delta Rule

The Delta rule is presented by Widrow and Hoff in 1960 as LMS (least mean square) learning procedure.

NN (Widrow's Adaline) is

$$\hat{y} = \sum_i w_i u_i + b$$

We want to adjust the weights to minimize the difference between the network output and the desired output. The error function

$$E(w_i) = \frac{1}{2} \sum_{k=1}^N [y(k) - \hat{y}(k)]^2$$

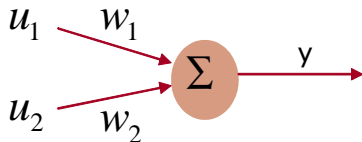
The idea is to make a change in the weight proportional to the negative of the derivative of the error as measured on the current data with respect to each weight: gradient descent.

Delta Rule

The weights are changed as

$$w_{new} = w_{old} - \eta \frac{\partial E}{\partial w}$$
$$\frac{\partial E}{\partial w_i} = \frac{\partial E}{\partial \hat{y}_i} \frac{\partial \hat{y}_i}{\partial w_i} = - [y(k) - \hat{y}(k)] u_i$$
$$w(k+1) = w(k) + \eta e(k) u_i$$

where $e(k) = y(k) - \hat{y}(k)$



$$w_1(k+1) = w_1(k) + \eta [y(k) - \hat{y}(k)] u_1$$

$$w_2(k+1) = w_2(k) + \eta [y(k) - \hat{y}(k)] u_2$$