

Review 1 (Chapters. 1, 2, 6, 3)**1. Basics**

- Comparison between human brain and von Neumann architecture
- Processing units
- Activation/output functions (threshold, linear-threshold, sigmoid)
- Network architecture (hidden nodes, feed-forward/recurrent nets, layered)
- Connection and weights
- Types of learning (supervised/unsupervised), Hebbian rule

2. Single Layer networks (Perceptron, Adaline, and the delta rule)

- Architecture
- Decision boundary and the problem of linear separability ($b + \sum_{i=1}^n x_i w_i = 0$)
- Hebbian nets ($\Delta w_i = x_i \cdot t$)
- Perceptron learning rule (only when $t \neq y : \Delta w_i = a \cdot x_i \cdot t$)
- Perceptron convergence theorem
- Delta learning rule in Adaline (driven by error: $\Delta w_i = a \cdot x_i \cdot (t - y_{in})$)
- Gradient descent approach in deriving delta learning rule
squared error: $E = (t - y_{in})^2$ or $E = \sum_{p=1}^P (t(p) - y_{in}(p))^2$
 $\Delta w_i \propto -\partial E / \partial w_i$

3. Backpropagation (BP) Networks

- Multi-layer feed-forward architecture with hidden nodes of non-linear and differentiable activation functions
- Motivation to have hidden nodes (representational power). Why non-linear?
- Feed forward computing
- BP learning
 - Training samples
 - Obtain errors at output layer (feed-forward phase): $d_k = (t_k - y_k) f'(y_{in_k})$
 - Obtain errors at hidden layer (error backpropagation phase): $d_j = d_{in_j} \cdot f'(y_{in_j})$
and $d_{in_j} = \sum_{k=1}^m d_k w_{jk}$
 - Learning procedure (batch and sequential modes)
 - In what sense BP learning generalizes delta rule of Adaline
 - Why BP learning works (gradient descent to minimize error): $\Delta w_{ij} = -a \cdot \partial E / \partial w_{ij}$
- Issues of practical concerns
 - Bias, error bound, training data, initial weights, number and size of hidden layers;
 - Learning rate (momentum, adaptive rate)
- Advantages and problems with BP learning
 - Powerful (general function approximator); easy to use; wide applicability; good generalization
 - Local minima; overfitting; parameters may be hard to determine; network paralysis; long learning time, hard to accommodate new samples (non-incremental learning)

4. Pattern Association and Associative memory (AM)

- Simple AM
Associative memory (AM) (content-addressable/associative recall; pattern correction/completion)
Network architecture: single layer or two layers of non-linear units
- Hebbian rule: $w_{ij} = \sum_{p=1}^P s_i(p) t_j(p)$
 - Correlation matrix: $W = \sum_{p=1}^P s^T(p) t(p)$ (assuming both s and t are row vectors).
 - Principal and cross-talk term: $s(k)W = \|s(k)\| + \sum_{p \neq k} s(k) s^T(p) t(p)$
- Delta rule: $\Delta w_{ij} = a \cdot (t_j - y_j) \cdot x_i$ or $\Delta w_{ij} = a \cdot (t_j - y_j) \cdot x_i \cdot f'(y - in_j)$ derived following gradient descent ($\Delta w_{ij} = -a \cdot \partial E / \partial w_{ij}$)
- Auto-associative memory: $t(p) = s(p)$, $p = 1, 2, \dots, P$.
Storage capacity: up to $n-1$ mutually orthogonal patterns of dimension n
- Iterative autoassociative memory
 - Motives (comparing with non-iterative recall)
 - Using the output of the current iteration as input of the next iteration (stop when a state repeats)
 - Dynamic system (stable states, attractors, genuine and spurious memories)
- Hopfield model for autoassociative memory
 - Network architecture (single-layer, fully connected, recurrent)
 - Weight matrix for Hopfield model (symmetric with zero diagonal elements)

$$w_{ij} = \begin{cases} \sum_{p=1}^P s_i(p) \cdot s_j(p) & \text{if } i \neq j \\ 0 & \text{otherwise} \end{cases}$$
 - Recall procedure (iterative until stabilized)
 - Stability of dynamic systems
 - Ideas of Lyapunov function/energy function (monotonically non-increasing and bounded from below)
 - Convergence of Hopfield AM: its an energy function
 - $E = -0.5 \sum_{i \neq j} \sum_j y_i y_j w_{ij} - \sum_i x_i y_i + \sum_i q_i y_i$
 - Storage capacity of Hopfield AM ($P \approx n / (2 \log_2 n)$).
- Bidirectional AM (BAM)
 - Architecture: two layers of non-linear units
 - Weight matrix: $W_{n \times m} = \sum_{p=1}^P s^T(p) t(p)$
 - Recall: bi-directional (from x to y and y to x); recurrent
 - Analysis
 - Energy function: $L = -XY^T = -\sum_{j=1}^m \sum_{i=1}^n x_i w_{ij} y_j$
 - Storage capacity: $P = O(\max(n, m))$