**(1) Epsilon-Greedy**

$$a_t = \begin{cases} \arg\max_i Q_{t-1}(i), & \text{with probability } 1 - \varepsilon, \\ \text{a random arm in } \{1, \ldots, K\}, & \text{with probability } \varepsilon, \end{cases}$$

$$N_t(a_t) = N_{t-1}(a_t) + 1, \quad Q_t(a_t) = Q_{t-1}(a_t) + \frac{1}{N_t(a_t)}\big(r_t - Q_{t-1}(a_t)\big).$$

**(2) UCB (Upper Confidence Bound)**

$$a_t = \arg\max_i \left[ Q_{t-1}(i) + c\sqrt{\frac{\ln t}{N_{t-1}(i)}} \right],$$

$$N_t(a_t) = N_{t-1}(a_t) + 1, \quad Q_t(a_t) = Q_{t-1}(a_t) + \frac{1}{N_t(a_t)}\big(r_t - Q_{t-1}(a_t)\big).$$

**(3) Softmax Action Selection**

$$P_t(i) = \frac{\exp\big(Q_{t-1}(i)/\tau\big)}{\sum_{j=1}^K \exp\big(Q_{t-1}(j)/\tau\big)}, \quad a_t \sim P_t(\cdot),$$

$$N_t(a_t) = N_{t-1}(a_t) + 1, \quad Q_t(a_t) = Q_{t-1}(a_t) + \frac{1}{N_t(a_t)}\big(r_t - Q_{t-1}(a_t)\big).$$

**(4) Thompson Sampling**

$$\theta_i \sim \text{Beta}\big(\alpha_i, \beta_i\big), \quad a_t = \arg\max_i \theta_i,$$

$$\alpha_{a_t} \leftarrow \alpha_{a_t} + r_t, \quad \beta_{a_t} \leftarrow \beta_{a_t} + (1 - r_t).$$

( $r_t \in \{0, 1\}$ Bernoulli )