

가치 함수(Q-table)의 업데이트는 ϵ -Greedy 정책과 Sample-average 방식을 결합하여 이루어진다. 가치 함수는 식(2)과 같이 갱신된다.

$$Q_{t+1}(a_{greedy}) = Q_t(a_{greedy}) + \frac{1}{N(a_{greedy})} [R - Q_t(a_{greedy})] \quad (2)$$

식(2)에서 $Q(a_{greedy})$ 는 행동 a_{greedy} 에서의 가치 함수를, $N(a_{greedy})$ 는 행동 a_{greedy} 를 선택한 누적 횟수를 나타낸다. 매 에피소드마다 확률 ϵ 로 탐험 (exploration)을, 확률 $(1-\epsilon)$ 로 Q-table에서의 최적의 행동 (활용, exploitation)을 선택하며, 탐험 비율(ϵ)은 에피소드가 진행될수록 점진적으로 감소시키는 방식을 사용하여, 초반의 다양한 탐색으로부터 후반으로 갈수록 학습된 최적 a_{greedy} 값을 더 많이 활용하게 된다. ϵ 의 감소 공식은 식(3)과 같다.

$$\epsilon \leftarrow \max(\beta\epsilon, \epsilon_{\min}) \quad (3)$$

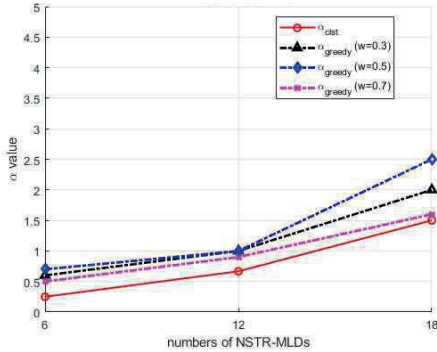
식(3)의 탐험 비율에 감소율(β)을 곱한 값과 탐험 비율 최솟값 중 최댓값으로 선택하여 총 70회의 에피소드를 반복 수행한 뒤, 학습된 Q-table에서 $Q(a_{greedy})$ 가 가장 높은 가치 함수를 갖는 a_{greedy} 값을 최적 파라미터로 선정하였다.

III. 모의실험

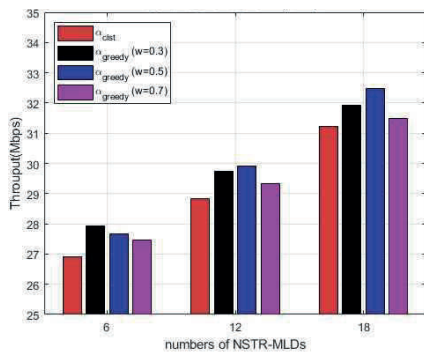
<표 1> 모의실험 환경

파라미터	값
모의실험 시간	1 sec
주파수 대역	2.4 GHz (HCL), 5 GHz (MDL)
다중 링크 수	2 개
단말 수 (HCL)	30 대
MCS	256 QAM (98 Mb/s)
MPDU	1000 bytes
대역폭	20 MHz
SIFS	18 μ sec
에피소드 수	70 회
ϵ	1.0 \rightarrow 0.01
β	0.98

<표 1>은 본 연구의 모의실험에 사용된 주요 파라미터를 나타낸다. 모의실험은 HCL에 연결된 전체 단말 (SLD와 NSTR-MLD) 수를 30으로 고정하고 SLD와 NSTR-MLD 수를 조정하면서 진행되었다. SLD 단말은 HCL에서만 동작한다고 가정하였다. 성능 평가 지표로는 HCL 링크의 평균 처리율과 단말 단위의 공정성 지수를 사용하였다.



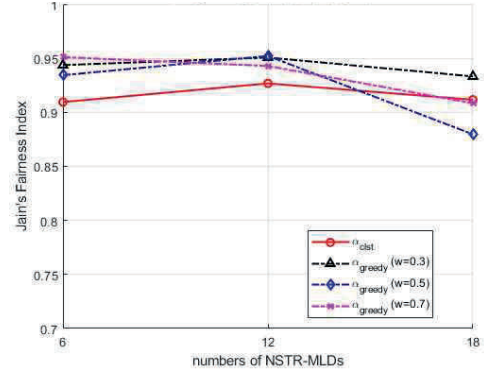
[그림 2] α_{clst} 와 α_{greedy} 값 변화 그래프



[그림 3] NSTR-MLD 단말 수 변화에 따른 처리율

[그림 2]는 NSTR-MLD 단말 수 변화에 따른 α_{clst} 와 α_{greedy} 값의 변화를 나타낸다. α_{clst} 값은 [1]에 주어진 것과 동일하게 NSTR-MLD 수를 SLD 수로 나눈 값으로 설정하였다. ϵ -Greedy 정책을 통해 학습된 α_{greedy} 값은 대체로 α_{clst} 값보다 큰 값을 가지며 가중치 w 에 따라 달라지는데, NSTR-MLD 수가 6개인 상황에서 w 가 커질수록 α_{clst} 대비 높은 값으로 도출되었다.

[그림 3]은 NSTR-MLD의 단말 수 변화에 따른 처리율을 보여준다. NSTR-MLD의 단말 수가 6 이고 $\alpha_{greedy}(w=0.3)$ 일 때 α_{clst} 대비 약 5% 높아진 처리율을 달성했다. 높은 α_{greedy} 값은 HCL과 MDL 링크가 백오프 과정 없이 동시 전송을 횟수가 많아지는 것을 의미한다. 이는 학습된 α_{greedy} 값이 적은 개수의 NSTR-MLD 환경에서 적극적인 연속 전송 전략을 선택하여 네트워크의 전반적 처리율을 증가시킨 결과이다.



[그림 4] Jain's Fairness Index 비교

[그림 4]는 NSTR-MLD의 단말 수 변화에 따른 공정성 지표를 보인다. HCL 처리율만을 고려할 경우 공정성의 악화를 초래할 수 있으므로 보상 함수에 공정성 지수를 반영하여 실험한 결과이다. 대체로 α_{clst} 보다 높은 값을 보여주지만 $\alpha_{greedy}(w=0.5)$ 일 때 단말 수가 늘어남에 따라 공정성이 다소 감소하는 것을 확인할 수 있다.

IV. 결론

본 연구는 Wi-Fi 7 다중 링크 동작 환경에서 CLST 기법의 핵심 파라미터인 α_{clst} 의 최적값을 얻기 위해 ϵ -Greedy 정책 기반의 Multi-Armed Bandit 학습 방법을 제안하였다. 모의실험을 수행한 결과, 학습된 α_{greedy} 값은 α_{clst} 방식 대비 HCL의 처리율과 공정성 지수를 모두 개선하였다. 본 연구의 시뮬레이션 결과는 에피소드 수와 시나리오 수의 제한으로 일부 한계가 있었으나, 추후 더욱 다양한 조건에서 추가적인 실험을 통해 더욱 개선된 결과를 얻을 수 있을 것으로 기대된다. 향후 연구에서는 ϵ -Greedy 외의 다양한 강화학습 알고리즘을 비교 및 평가하고, 실시간 네트워크 상태 변화에 따라 α_{greedy} 값을 동적으로 최적화하는 방법을 추가로 개발하여, 본 연구의 실질적인 적용성을 더욱 높이고자 한다.

ACKNOWLEDGMENT

이 논문은 정부(과학기술정보통신부)의 재원으로 정보통신기획평가원-혁신사업ICT핵심인재양성 지원을 받아 수행된 연구임 (IITP-2024-00436744)

참 고 문 헌

- [1] Kwon, Lam, and Eun-Chan Park. "Contention-Less Multi-Link Synchronous Transmission for Throughput Enhancement and Heterogeneous Fairness in Wi-Fi 7." *Sensors* 24.11 (2024): 3642.
- [2] Luong, Nguyen Cong, et al. "Applications of deep reinforcement learning in communications and networking: A survey." *IEEE communications surveys & tutorials* 21.4 (2019): 3133-3174.