

we see a similar outcome for average rewards. Optimistic and Realistic greedy eventually converge, while VCB performs better due to ~~optimized~~ better exploration and ceasing exploration when a good estimate has been found.

2. In Non-stationary case.

Optimal Action: VCB performance ~~keeps~~ is bad because it stops exploring actions after ^(~2000) some steps have passed.

It starts acting greedily and its estimates do not get updated with ~~as~~ a constant step-size.

Optimistic initially performs well. This is due to optimal action being picked more often. However because of the initial bias, there will be significant variance in ~~the~~ estimates of actions other than the supposed "optimal" estimate if got early in the simulation.

Realistic eventually performs better because it keeps exploring and ~~thus~~ the samples the ~~new~~ distributions is updated throughout.

Argo Rewards : Optimistic performs best

due to initial bias and selection of the optimal value.

UCB performs better than Realistic because of initially exploring