1- montecarlo would be better in the case our prior state-values are not much more reliable because of some changes. So montecarlo would understand the change more rapidly and avoid bias.

2- a) Because in Q-lerng we are estimaty $Q^*$; the optimal action-value and we allecate $a's$ argmax $Q(s,a)$ for next iteration, Independent of the policy beg followed.

2-b) No. Because Q-leay is still off-policy and for behairior policy could select any policy ($\varepsilon$-greedy) while SARSA is on-policy and greed selection cause it not have exploration so weight updates will be different

3) a) No I do not think So.

No I don't think other differant value of $\alpha$ would have better result. either larger and smaller $\alpha$ have some disadvantages, for example smaller $\alpha$ will take long time for convergence so although it may reach a closer approximation of optimal value but it takes too long time. on the other hand, large $\alpha$ will cause noisy convergence and some distance from the
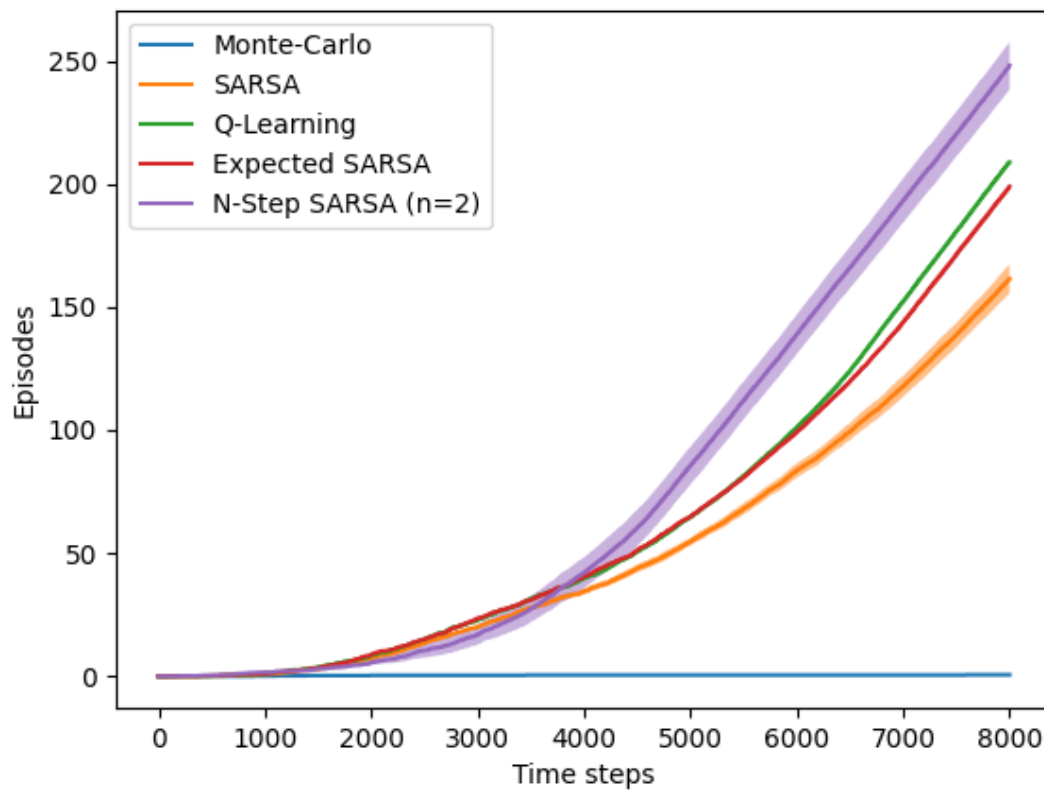
3-b) the larger $\alpha$ the more dependancy to next value function So in this case when $\alpha$ is lang and the initial values are selected less than optimal value at initial steps our estimation will more under optimal value (under shooting) and then will be corrected. ~~of~~ ~~~~

3-6) Because we are using n-step we want to avoid finish early for random walk. If we have small problem, the expected value of number of states of traversely will be smaller than n, making the performance less fair for larger n so for very small problem smaller n would be better
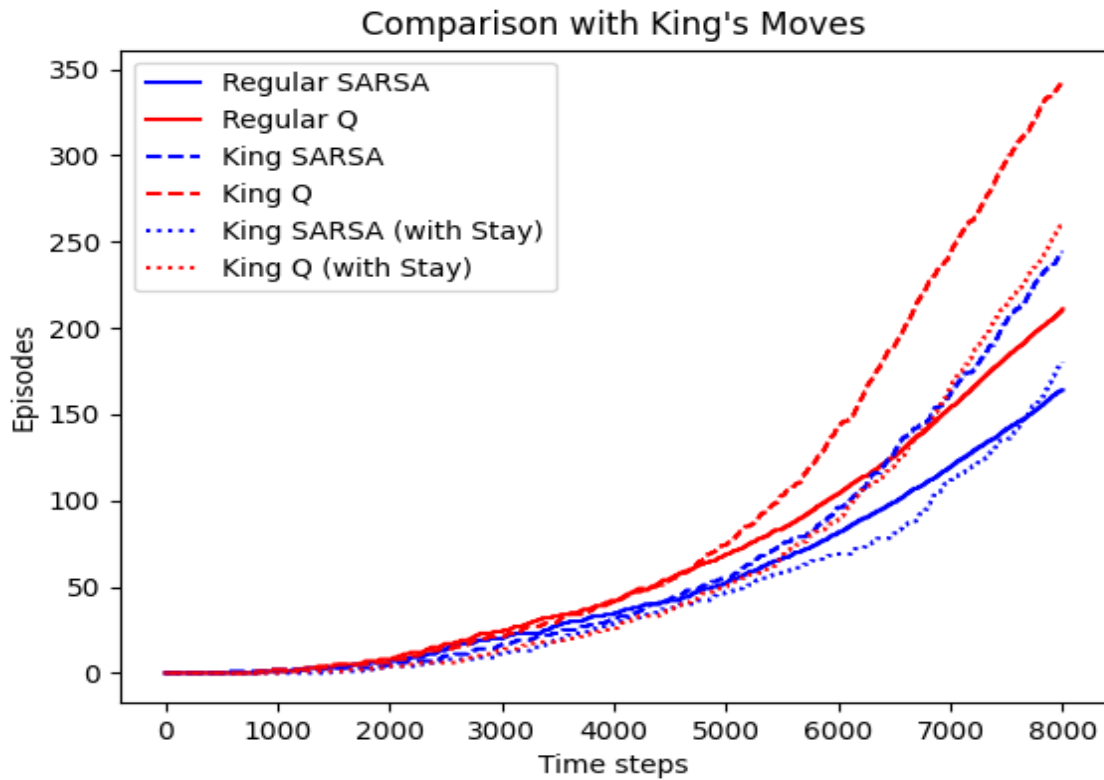
**The algorithm related to Q4 and Q5 are in Algorithms file. Right now, 4a is uncommented in the main function at the bottom of the file.**
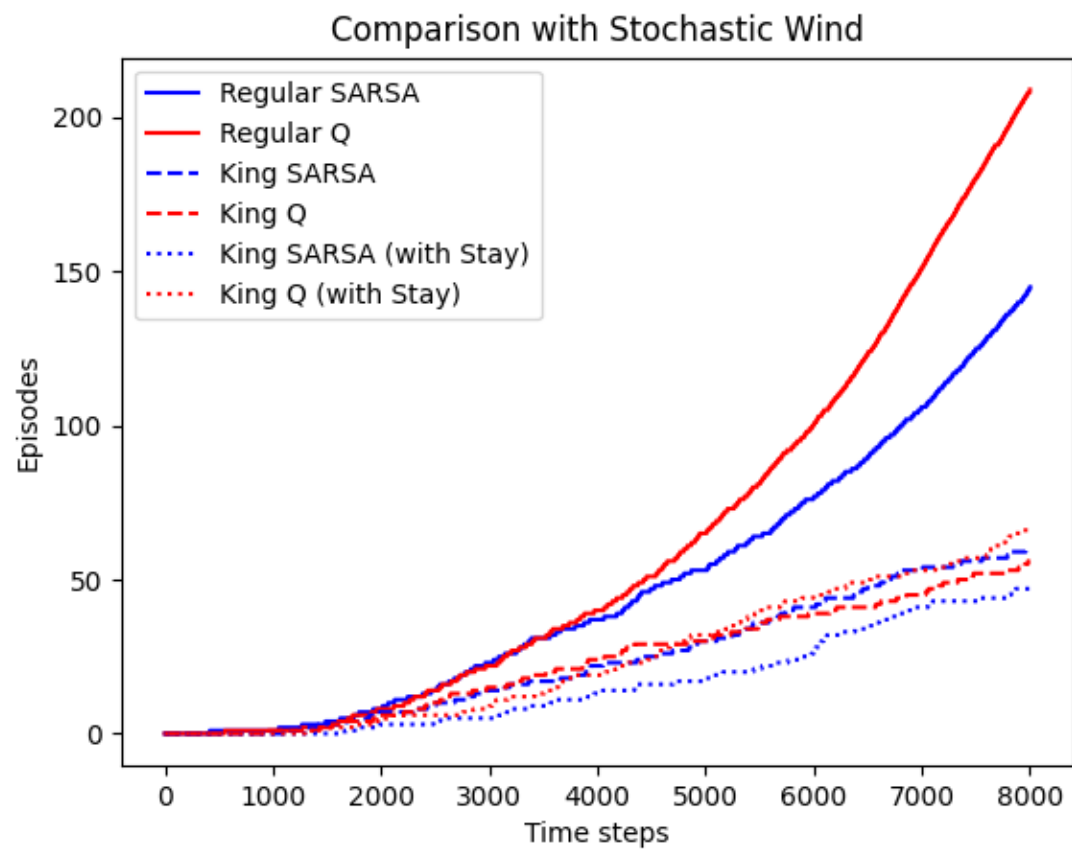
# Q4)

**4-b-a)** I used n = 2 for N-Step SARSA because it gave me better performance

**4-c)** Having more actions did give us better performance, as we could continue to move straight right even through weak wind by using the down-right action. Adding the stay action reduced performance, as I believe this would be a trap action to take, which means the agent wastes time learning that it is a bad action.
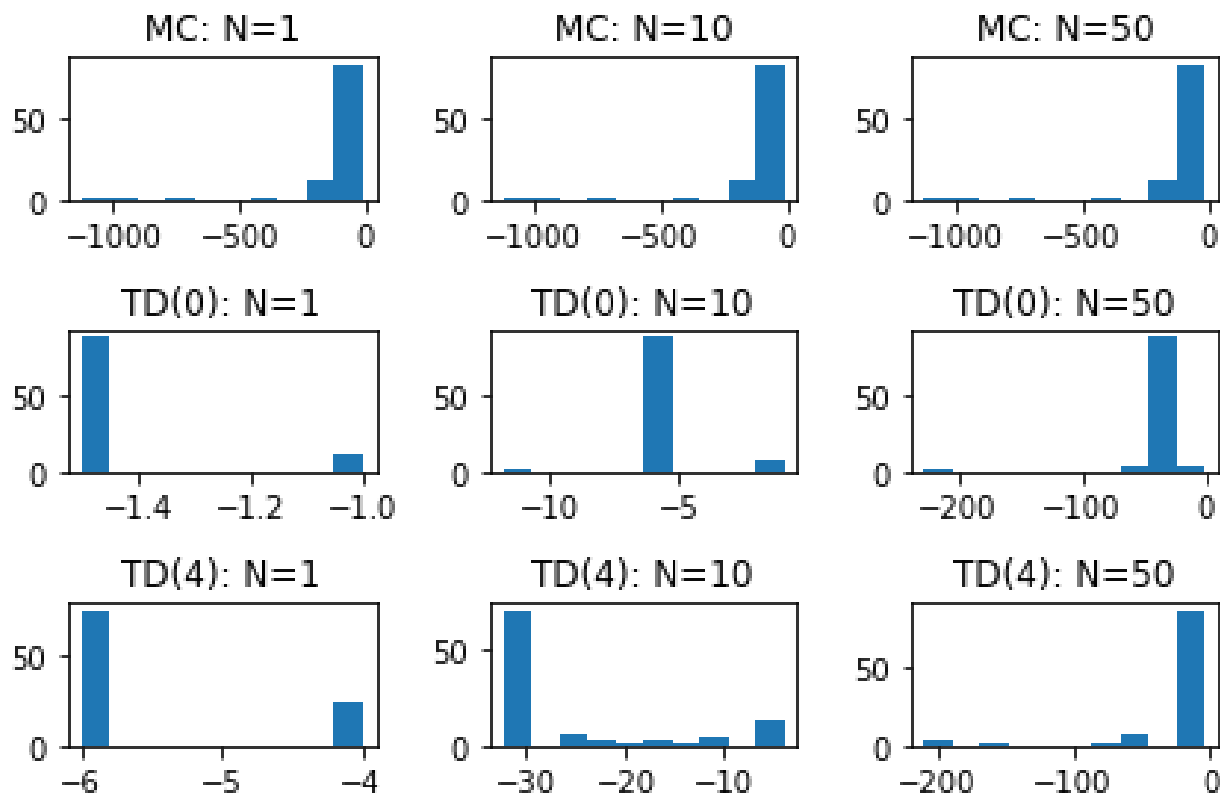


Comparison with King's Moves

**4-d)** Stochastic wind causes the problem to be harder to solve

Comparison with Stochastic Wind

**Q 5)**

**5-a)**

**5-b)**

There seems to be high variance in the Monte-Carlo episodes, but no bias, as the plots do not change with more training episodes. On the other hand, TD methods have incredibly small variance, but high bias, as they continue to change with more training episodes