

Reinforcement Learning

2nd Programming Assignment

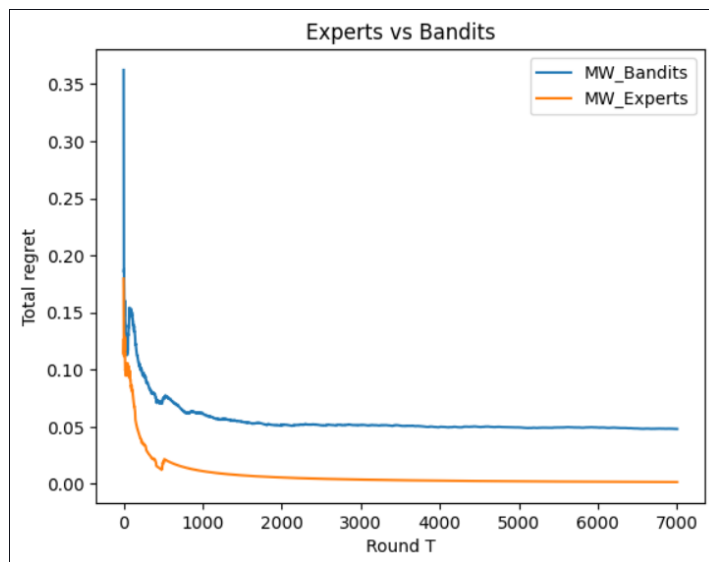
ΚΩΝΣΤΑΝΤΙΝΟΣ ΝΙΚΟΛΟΣ
AM:2019030096

Περιγραφή της άσκησης:

Ζητείται η υλοποίηση του αλγορίθμου Multiply Weights σε περιβάλλον Experts και Bandits. Βασική διαφορά των 2 αποτελεί η έλλειψη γνώσης για όλους τους experts. Στο εχθρικό περιβάλλον Bandits γνωρίζουμε το loss μόνο του expert που επιλέχθηκε σε κάθε γύρο επομένως αλλάζουμε μόνο εκείνη την τιμή καθώς και μόνο το δικό του βάρος. Αποτέλεσμα αυτού η αργότερη 'μάθηση' σε περιβάλλον bandits σε αντίθεση με τους Experts.

Γνωρίζουμε ότι το Multiply Weights Algorithm σε Experts περιβάλλον αναμένεται να έχει $\text{Regret}T = \leq 2\sqrt{T\ln k}$.

Ενώ σε Bandits αναμένεται να έχει $\text{Regret}T = \sqrt{kT\ln k} = 5.47\sqrt{T\ln k} = 2.73 * 2\sqrt{T\ln k}$ γεγονός που μπορεί να επιβεβαιωθεί από το παρακάτω διάγραμμα



Ακόμη ένα μέτρο σύγκρισης που μας ενδιαφέρει είναι πως στο Experts παρατηρούνται 4820 σωστές προβλέψεις σε αντίθεση με τις 3711

