maze learning: optimal pseudorewards during some episodes Q-learning use pseudorewards on episode with prob 0.1 prob 0.25 prob 0.5 prob 0.9 steps episode