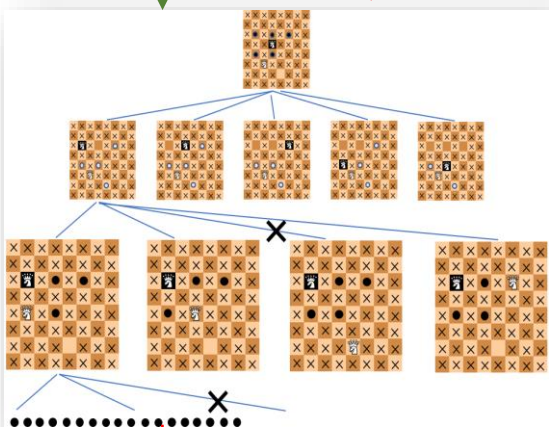


奖励估值

初始策略

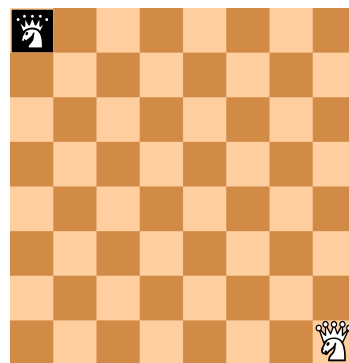
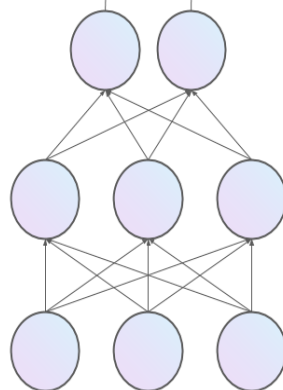
策略分布

奖励值



提高的策略

Cross  
Entropy



Self Play

游戏结果

MSE