

Value iteration implementation for the Gambler's Problem

در این سوال قصد داریم در قالب حل مسئله قمارباز (Gambler Problem) با الگوریتم value iteration آشنا شویم.

مسئله مطابق با فرض های صورت تمرین تعریف می شود. مطابق با الگوریتم value iteration در پیاده سازی، state ها، دارایی هستند. همچنین action ها مقدار سرمایه ای است که روی آن شرط بندی می شود. برای ارزش state ها یک آرایه با مقادیر اولیه صفر تعریف می شود. همچنین در نهایت، بهترین action در هر state را در یک آرایه (policy) ذخیره می کنیم. همچنین آرایه ای که مطابق با پاداش هر state است در تمامی اندیس ها صفر بوده و تنها در state آخر (پایان بازی) ۱ می باشد. برای پیاده سازی الگوریتم مطابق با تصویر ۱-۲ عمل می کنیم. به این منظور که در یک حلقه بی نهایت یک Δ با مقدار اولیه صفر تعریف می کنیم و سپس بعد از به روز کردن ارزش هر state مقدار آن را با اختلاف ارزش قبلی و بعدی state جایگزین می کنیم. (در صورتی که این مقدار از مقدار سابق Δ کوچکتر باشد، جایگزینی صورت نمی گیرد.) در حقیقت Δ بزرگترین تغییر ارزش state ها را در هر بار به روزرسانی تمامی ارزش ها نشان می دهد. برای به روز کردن ارزش ها مطابق با رابطه تصویر ۱-۲، عمل می کنیم. نکته قابل توجه این است که از یک state خاص نمی توانیم به تمامی state های دیگر برویم. زیرا در صورتی که کمتر از 50\$ داشته باشیم نمی توانیم روی مقداری بیشتر از سرمایه خود شرط بندی کنیم و در صورتی که بیشتر از 50\$ نیز داشته باشیم، معقولانه است که نهایتاً طوری شرط بندی کنیم که به 100\$ برسیم. در نتیجه state های بعدی را فقط state هایی در نظر می گیریم که در این دو شرط صدق کند. برای به روز کردن ارزش ها تابع $\text{Best_NextState}()$ نوشته شده است. در انتها نیز برای انتخاب بهترین action در هر state یکبار $\text{Best_NextState}()$ را فراخوانی می کنیم. این الگوریتم به طور کلی در تابع $\text{Value_Iteration}()$ پیاده سازی شده است.

Initialize array V arbitrarily (e.g., $V(s) = 0$ for all $s \in \mathcal{S}^+$)

Repeat

$\Delta \leftarrow 0$

For each $s \in \mathcal{S}$:

$v \leftarrow V(s)$

$V(s) \leftarrow \max_a \sum_{s',r} p(s',r|s,a)[r + \gamma V(s')]$

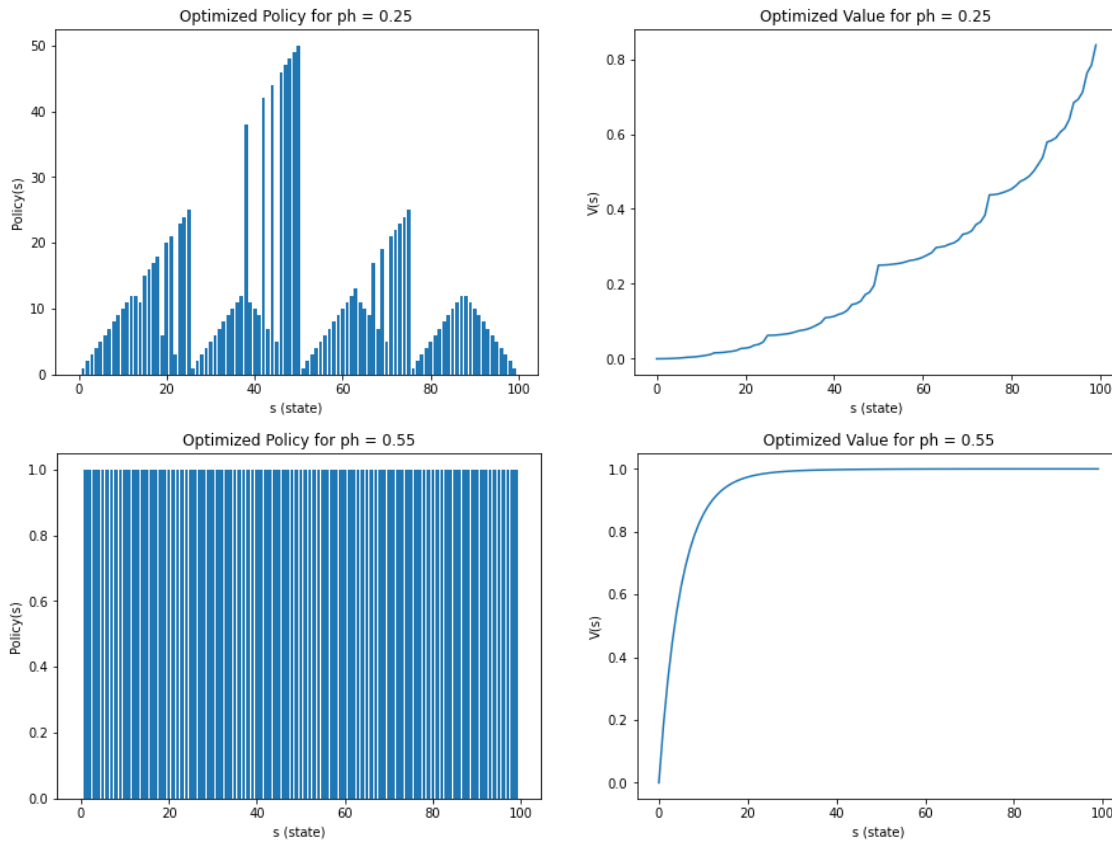
$\Delta \leftarrow \max(\Delta, |v - V(s)|)$

until $\Delta < \theta$ (a small positive number)

Output a deterministic policy, π , such that

$\pi(s) = \arg \max_a \sum_{s',r} p(s',r|s,a)[r + \gamma V(s')]$

تصویر ۲-۲ نمودار های ارزش و policy برای دو احتمال مختلف $p_h = 0.25$ و $p_h = 0.55$ نشان می دهد.



تصویر ۲-۲: نمودار های ارزش و policy در مسئله قمارباز برای دو احتمال مختلف

همانطور که در تصاویر بالا دیده می شود به نکات زیر می رسیم:

- ۱- هر چه به state های انتهایی نزدیک تر می شویم، ارزش آن بیشتر می شود.
- ۲- اگر احتمال شیر آمدن بزرگتر از 0.5 باشد، می دانیم که در هر صورت (اگر چه کند باشد) ولی بالاخره برنده بازی خواهیم شد و بنابراین همیشه روی 1\$ شرط بندی می کنیم. در نتیجه هنگامی که در state های بالاتر هستیم در صورت باخت ارزش state چندان تغییری نمی کند. به بیان دیگر چون state های مجاور در صورت برد یا باخت عوض می شوند (چون روی 1\$ شرط بندی کرده ایم، تغییری در ارزش های state های بالاتر احساس نمی شود).
- ۳- اگر احتمال شیر آمدن کوچکتر از 0.5 باشد، می دانیم که شانس برنده شدن در هر بار کم است؛ در نتیجه روی مقدار زیادی شرط بندی می کنیم. یعنی اگر تعداد دفعات زیادی شرط بندی کنیم می دانیم در تعداد کمی از آنها برنده هستیم. پس در بعضی مواقع مثلاً state های 25 و 50 روی کل پول شرط بندی می کنیم تا زودتر برنده یا بازنده شویم. (اگر تعداد دفعات شرط بندی

زیاد باشد در هر صورت تعداد باخت ها از برد ها بیشتر است. اگر تعداد دفعات کم باشد، ممکن است با یک شیر آمدن بازی را برنده شویم.