

هوش مصنوعی

جزوه پنجم

بازی ها:

در یادداشت اول، ما در مورد مسائل جستجو و نحوه حل آنها به طور کارآمد و بهینه صحبت کردیم - با استفاده از الگوریتم های قدرتمند جستجوی تعمیم یافته، عوامل ما می توانند بهترین طرح ممکن را تعیین کنند و سپس به سادگی آن را برای رسیدن به یک هدف اجرا کنند. حال، بیایید دنده ها را تغییر دهیم و سناریوهایی را در نظر بگیریم که در آن عوامل ما یک یا چند دشمن دارند که سعی می کنند آنها را از رسیدن به اهدافشان باز دارند. عوامل ما دیگر نمی توانند الگوریتم های جستجویی را که قبلاً یاد گرفته ایم برای تدوین یک برنامه اجرا کنند، زیرا ما معمولاً نمی دانیم که دشمنان ما چگونه علیه ما برنامه ریزی می کنند و به اقدامات ما پاسخ می دهند. در عوض، ما نیاز به اجرای یک کلاس جدید از الگوریتم ها داریم که راه حل هایی را برای مسائل جستجوی خصمانه ارائه می دهند که معمولاً به عنوان بازی شناخته می شوند.

انواع مختلفی از بازی ها وجود دارد. بازی ها می توانند کنش هایی با نتایج قطعی یا تصادفی (احتمالی) داشته باشند، یا می توانند تعداد بازیکنان متغیری داشته باشند و یا ممکن است مجموع صفر باشند یا نباشند. اولین دسته از بازی هایی که به آن خواهیم پرداخت، بازی های مجموع صفر قطعی هستند، بازی هایی که در آنها اقدامات قطعی هستند و سود ما مستقیماً معادل باخت حریف است و بالعکس. ساده ترین راه برای فکر کردن در مورد چنین بازی هایی این است که با یک مقدار متغیر تعریف شوند، که یک تیم یا عامل سعی می کند آن را به حداکثر برساند و تیم یا عامل مقابل سعی می کند آن را به حداقل برساند و آنها را در رقابت مستقیم قرار دهد. در یک من، این متغیر امتیاز شما است که سعی می کنید با خوردن سریع و کارآمد گلوله ها آن را به حداکثر برسانید در حالی که ارواح سعی می کنند ابتدا شما را بخورند. بسیاری از بازی های خانگی رایج نیز در این دسته از بازی ها قرار می گیرند:

- چکرز^۱: اولین بازیکن کامپیوتری چکرز در سال ۱۹۵۰ ساخته شد. از آن زمان، چکرز به یک بازی حل شده تبدیل شده است، به این معنی که هر موقعیتی را می توان به عنوان یک برد، باخت یا مساوی برای هر دو طرف ارزیابی کرد، زیرا هر دو بازیکن به طور بهینه عمل می کنند.
- شطرنج: در سال ۱۹۹۷، دیپ بلو^۲ اولین عامل کامپیوتری شد که قهرمان شطرنج انسانی گری کاسپاروف را در یک مسابقه شش بازی ای شکست داد. دیپ بلو برای استفاده از روش های بسیار پیچیده برای ارزیابی بیش از ۲۰۰ میلیون موقعیت در ثانیه ساخته شده است. برنامه های فعلی حتی بهترند، هرچند کمتر تاریخی اند.

¹ adversaries

² adversarial search problems

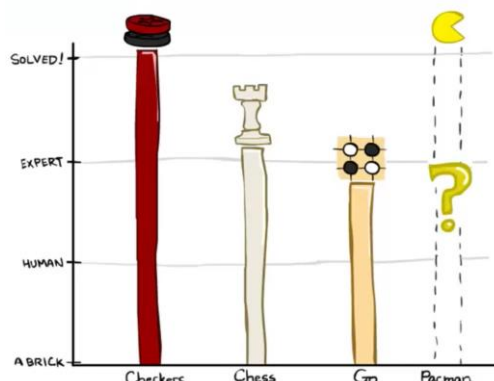
³ zero-sum

⁴ deterministic zero-sum games

⁵ Checkers

⁶ Deep Blue

- گو⁷: فضای جستجو برای این بازی بسیار بزرگتر از شطرنج است، و بنابراین اکثر آنها باور نداشتند که عاملان کامپیوتری هرگز قهرمانان جهان را برای چندین سال آینده شکست دهند. با این حال، AlphaGo، توسعه یافته توسط گوگل، قهرمان Go Lee Sodol را در مارس ۲۰۱۶، ۴ بر ۱ شکست داد.



همه عوامل قهرمان جهان که در بالا گفتیم، حداقل تا حدی، از تکنیک‌های جستجوی خصمانه استفاده می‌کنند که می‌خواهیم درباره آنها صحبت کنیم. برخلاف جستجوی معمولی که یک طرح جامع را برمی‌گرداند، جستجوی خصمانه یک استراتژی یا خط مشی را برمی‌گرداند که با توجه به پیکربندی‌های عامل(های) ما و مخالفان آنها، به سادگی بهترین حرکت ممکن را توصیه می‌کند. به زودی خواهیم دید که چنین الگوریتم‌هایی دارای ویژگی زیبایی هستند که از طریق محاسبات باعث ایجاد رفتار می‌شوند - محاسباتی که ما اجرا می‌کنیم از نظر مفهومی نسبتاً ساده و به طور گسترده قابل تعمیم است، با این حال به طور ذاتی یک همکاری بین عوامل در همان تیم و همچنین "تفکر" در مورد عوامل متخاصم ایجاد می‌کنند.

فرمول استاندارد بازی شامل تعاریف زیر است:

- حالت شروع: S_0
- بازیکنان: بازیکن (بازیکنان) مشخص می‌کنند که نوبت چه کسی است
- کنش‌ها: کنش(کنش‌ها) ی قابل اجرا برای هر بازیکن
- نتیجه مدل انتقال (s, a)
- بررسی پایانه: $\text{Terminal-test}(s)$
- مقادیر پایانه: $\text{Utility}(s, \text{player})$

مینیماکس (Minimax):

[استراتژی «حداقل-حداکثر» یا «کمین‌بیش»، یک استراتژی محتاطانه است که حداکثر ضرر را به حداقل می‌رساند. این استراتژی برای یک تصمیم گیرنده محافظه کار مفید است. اساساً این روش برای فرار از باخت و حداقل کردن ضررها صورت گرفته و مخصوص افرادی است که که نمی‌خواهند تصمیم اشتباه بگیرند.

⁷ Go

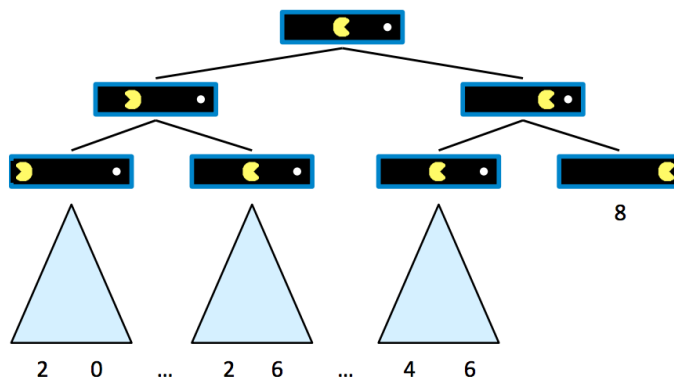
این الگوریتم تصمیم مینیماکس را از حالت فعلی محاسبه میکند و مقادیر مینیماکس مربوط به هر حالت پسین را بطور بازگشتی با پیاده سازی مستقیم معادلات تعریف شده، محاسبه میکند. این محاسبات بازگشتی به طرف پایین تا برگهای درخت پیش میرود و سپس در برگشت از حالت بازگشتی، مقادیر مینیماکس هر گره را اختصاص میدهد.

الگوریتم مینیماکس اکتشاف عمقی کاملی را روی درخت بازی انجام میدهد. اگر حداکثر عمق درخت m باشد و در هر نقطه b حرکت معتبر وجود داشته باشد، آنگاه پیچیدگی زمانی الگوریتم مینیماکس برابر است با $O(b^m)$. برای الگوریتمی که تمام فعالیت ها را همزمان تولید میکند، پیچیدگی فضا (حافظه) برابر با $O(bm)$ ، یا برای الگوریتمی که هر بار یک فعالیت را تولید میکند، برابر با $O(m)$ است. البته برای بازی های واقعی هزینه زمان کاملاً غیرعملی است. اما این الگوریتم به عنوان مبنایی برای تحلیل ریاضی بازی ها و برای بسیاری از الگوریتم های عملی است.⁸

اولین الگوریتم بازی مجموع صفر که در نظر خواهیم گرفت، مینیماکس است، که با این فرض انگیزشی اجرا می شود که حریفی که با آن روبرو هستیم رفتار بهینه ای دارد و همیشه حرکتی را که برای ما بدترین حالت است را انجام می دهد. برای معرفی این الگوریتم، ابتدا باید مفهوم مطیوبیت های پایانه⁹ و مقدار حالت¹⁰ را شرح دهیم. مقدار یک حالت امتیاز بهینه ای است که توسط عاملی که آن حالت را کنترل می کند به دست می آید. برای درک معنای این موضوع، صفحه بازی پک من ساده زیر را مشاهده کنید:



فرض کنید که پک من با ۱۰ امتیاز شروع می کند و در هر حرکت ۱ امتیاز از دست می دهد تا زمانی که گلوله را بخورد، در این مرحله بازی به حالت پایانی می رسد و بازی به پایان می رسد. می توانیم به شکل زیر شروع به ساخت درخت بازی برای این بازی کنیم، که در آن فرزندان یک حالت، وضعیت های جانشین هستند، درست مانند درخت های جستجو برای مسائل جستجوی عادی:



از این درخت مشخص است که اگر پکمن مستقیماً به سمت گلوله برود، بازی را با امتیاز ۸ به پایان می رساند، در حالی که اگر در هر نقطه عقب نشینی کند، در نهایت با مقداری با امتیاز کمتر به پایان می رسد. اکنون که یک درخت بازی با چندین حالت پایانه و واسطه ایجاد کرده ایم، آماده هستیم تا معنای ارزش هر یک از این حالت ها را شرح دهیم.

ارزش یک حالت به عنوان بهترین نتیجه ممکن (مطلوب) تعریف می شود که یک عامل می تواند از آن حالت به دست آورد. ما مفهوم مطلوبیت را بعداً به طور ملموس تر شرح می دهیم، اما در حال حاضر کافی است به سود یک عامل به عنوان امتیاز یا تعداد امتیازهایی که به دست

⁸ terminal utilities

⁹ state value

¹ successor states

می‌آورد فکر کنیم. مقدار یک حالت پایانه، که مطلوبیت پایانه^۱ نامیده می‌شود، همیشه مقداری از مقادیر شناخته شده قطعی و یک ویژگی ذاتی از بازی است. در مثال پک من ما، مقدار سمت راست ترین حالت پایانه به سادگی ۸ است، امتیازی که پک من با رفتن مستقیم به سمت گلوله به دست می‌آورد. همچنین در این مثال مقدار یک حالت غیر پایانی به عنوان حداکثر مقادیر فرزندان آن تعریف شده است. با تعریف $V(s)$ به عنوان تابعی که مقدار یک حالت s را تعریف می‌کند، می‌توانیم مبحث فوق را خلاصه کنیم:

$$\forall \text{ non-terminal states, } V(s) = \max_{s' \in \text{successors}(s)} V(s')$$

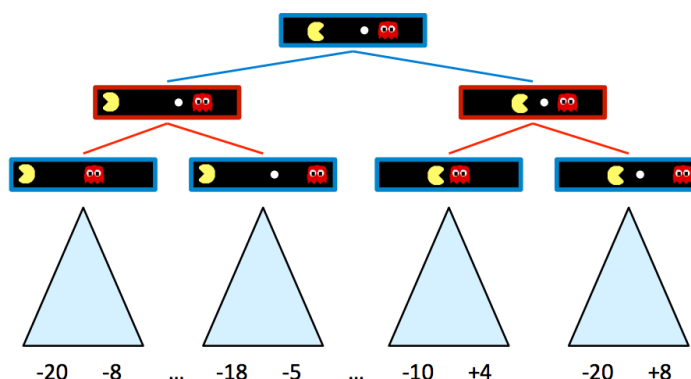
$$\forall \text{ terminal states, } V(s) = \text{known}$$

این یک قاعده بازگشتی بسیار ساده را تنظیم می‌کند، که طبق آن منطقی است که مقدار فرزند راست مستقیم گره ریشه ۸ باشد و فرزند سمت چپ مستقیم گره ریشه ۶ باشد، زیرا اینها حداکثر امتیازات احتمالی هستند که عامل می‌تواند به دست آورد، اگر به ترتیب از حالت شروع به سمت راست یا چپ حرکت کند. نتیجه این است که با اجرای چنین محاسباتی، یک عامل می‌تواند بفهمد که حرکت به سمت راست بهینه است، زیرا فرزند سمت راست ارزش بیشتری نسبت به فرزند سمت چپ از حالت شروع دارد.

حالا بیایید یک صفحه بازی جدید با یک روح خصمانه معرفی کنیم که می‌خواهد جلوی پک من را از خوردن گلوله بگیرد.

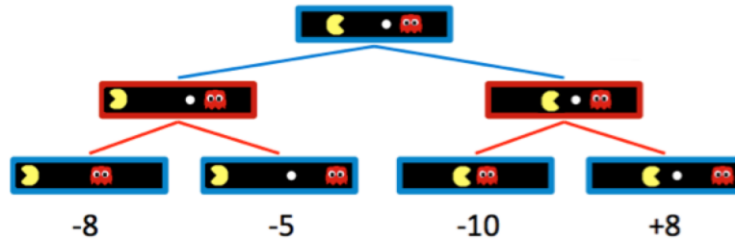


طبق قوانین بازی دو عامل به نوبت حرکت می‌کنند و به درخت بازی ای منتهی می‌شوند که در آن دو عامل لایه های درخت را که "تحت کنترل دارند" عوض می‌کنند. کنترل یک عامل بر روی یک گره، یعنی گره با حالتی مطابقت دارد که در آن نوبت آن عامل است، و بنابراین این فرصت برای اوست که در مورد یک عمل تصمیم بگیرد و وضعیت بازی را بر اساس آن تغییر دهد. در اینجا درخت بازی حاصل از صفحه جدید بازی دو عاملی بالا را داریم:



گره های آبی مربوط به گره های تحت کنترل پک من است و در آنها پک من می‌تواند تصمیم بگیرد که چه اقدامی انجام دهد، در حالی که گره های قرمز مربوط به گره های تحت کنترل روح است. توجه کنید که همه فرزندان گره های تحت کنترل روح، گره هایی هستند که روح از حالت قبلی خود در والدش به چپ یا راست حرکت کرده است، و همینطور این روند برای گره های تحت کنترل پک من نیز رخ داده است. برای سادگی، بیایید این درخت بازی را به درخت عمق ۲ تقلیل بدهیم و مقادیر جعلی را به حالت های پایانه ها به صورت زیر اختصاص دهیم:

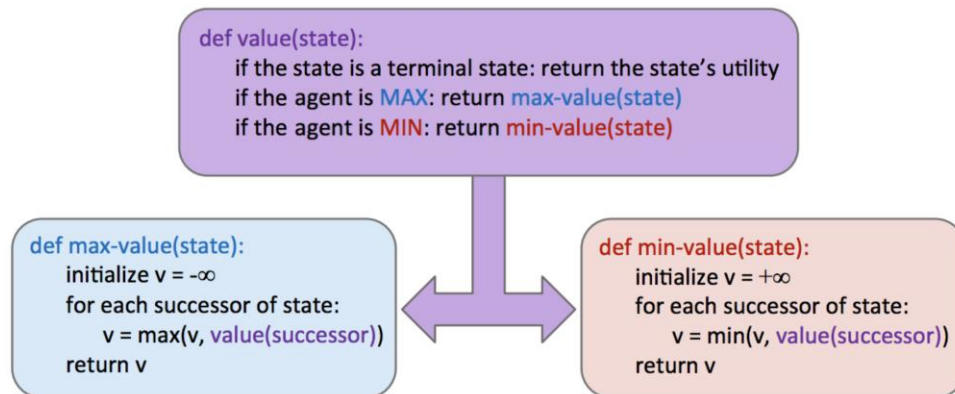
¹ terminal utility



طبیعتاً، افزودن گره‌های تحت کنترل روح، حرکتی را که پک من معتقد است بهینه است را، تغییر می‌دهد و حرکت بهینه جدید با الگوریتم مینیماکس تعیین می‌شود. الگوریتم مینیماکس به جای به حداکثر رساندن سودمندی بر روی فرزندان در هر لایه از درخت، فقط فرزندان گره‌های کنترل شده توسط پک من را به حداکثر می‌رساند، در حالی که فرزندان گره‌های تحت کنترل روح را به حداقل می‌رساند. از این رو، دو گره روح بالا به ترتیب دارای مقادیر $\min(-10, +8) = -10$ و $\min(-8, -5) = -8$ هستند. به همین ترتیب، گره ریشه نحن کنترل پک من دارای مقدار $\max(-8, -10) = -8$ است. از آنجایی که پکمن می‌خواهد امتیاز خود را به حداکثر برساند، به سمت چپ می‌رود و امتیاز ۸- را می‌گیرد به جای اینکه به دنبال گلوله برود و امتیاز ۱۰- را بگیرد. این نمونه بارز تغییر رفتار از طریق محاسبات است - اگرچه پکمن می‌خواهد امتیاز ۸+ را دریافت کند، ولی اگر در نهایت در سمت راست‌ترین حالت فرزند قرار بگیرد، از طریق مینیماکس «می‌داند» که یک روح با عملکرد بهینه به او اجازه نمی‌دهد آن امتیاز را کسب کند. به منظور انجام بهینه، پک من مجبور می‌شود شرط‌های خود را رعایت کند و به طور غیرمنتظره از گلوله دور شود تا کاهش امتیاز خود را به حداقل برساند. ما می‌توانیم روشی که مینیماکس مقادیر را به حالت‌ها اختصاص می‌دهد به صورت زیر خلاصه کنیم:

$$\begin{aligned} \forall \text{agent-controlled states, } V(s) &= \max_{s' \in \text{successors}(s)} V(s') \\ \forall \text{opponent-controlled states, } V(s) &= \min_{s' \in \text{successors}(s)} V(s') \\ \forall \text{terminal states, } V(s) &= \text{known} \end{aligned}$$

در پیاده‌سازی، مینیماکس مانند جستجوی اول عمق عمل می‌کند و مقادیر گره‌ها را به همان ترتیبی که DFS انجام می‌دهد محاسبه می‌کند، از سمت چپ‌ترین گره شروع می‌کند و به طور مکرر به سمت راست حرکت می‌کند. به طور دقیق‌تر، پیمایش پس‌ترتیب از درخت بازی را انجام می‌دهد. شبه‌کد به دست آمده برای مینیماکس هم زیبا و هم از نظر شهودی ساده است و در زیر ارائه شده است. توجه داشته باشید که مینیماکس یک عمل را که مربوط به گره ریشه است را به فرزندی که مقدار خود را از آن گرفته است، برمی‌گرداند.



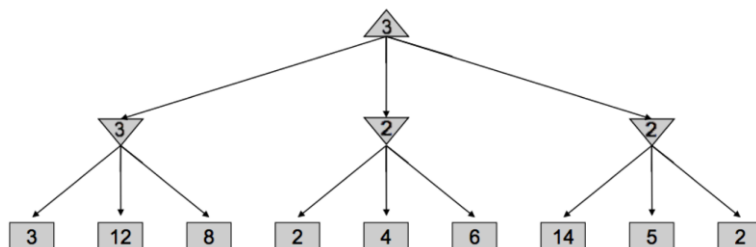
هرس آلفا-بتا:

[هرس آلفا-بتا (alpha-beta pruning) روشی است برای یافتن جواب بهینه‌ی الگوریتم مینیماکس که از جست‌وجو در برخی فضاهایی که انتخاب نمی‌شوند، جلوگیری می‌کند. به این ترتیب جست‌وجو در درخت کمینه-بیشینه تا عمق مشخصی در زمان کمتری انجام می‌شود.

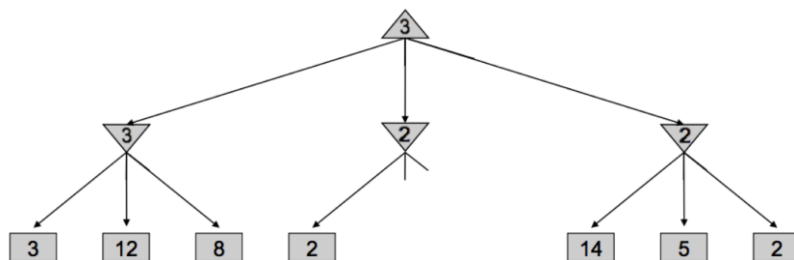
این الگوریتمی کارایی الگوریتم درخت مینیماکس (درخت کمینه بیشینه یا درخت بازی) را بهبود می‌بخشد. با استفاده از هرس آلفا-بتا، بخش‌هایی از درخت کمینه بیشینه که پیمایششان بی‌تأثیر است پیمایش نمی‌شوند و به این ترتیب پیمایش درخت کمینه بیشینه تا یک عمق مشخص در زمانی کم‌تر صورت می‌گیرد.]

مینیماکس تقریباً کامل به نظر می‌رسد چراکه ساده، بهینه و شهودی است. با این حال، اجرای آن بسیار شبیه به جستجوی اول عمق است و پیچیدگی زمانی اش با آن یکسان است، یعنی $O(b^m)$. با یادآوری اینکه b عامل انشعاب و m عمق تقریبی درختی است که گره‌های پایانه را می‌توان در آن یافت، این الگوریتم زمان اجرای بسیار زیادی را برای بسیاری از بازی‌ها به همراه دارد. برای مثال، شطرنج با ضریب انشعاب $b \approx 35$ و عمق درخت $m \approx 100$ است پس پیچیدگی زمانی اش خیلی زیاد می‌شود. برای کاهش این مقدار، مینیماکس یک بهینه‌سازی دارد که هرس آلفا-بتا است.

از نظر مفهومی، هرس آلفا-بتا به این صورت است: اگر می‌خواهید ارزش گره n را با نگاه کردن به جانشین‌های آن تعیین کنید، به محض اینکه متوجه شدید که مقدار n می‌تواند در بهترین حالت با مقدار بهینه والد n برابری کند، دیگر به دنبال آن نگردید. با یک مثال معنای این جمله را شرح می‌دهیم. درخت بازی زیر را با گره‌های مربعی مربوط به حالت‌های پایانی، مثلث‌های رو به پایین مربوط به گره‌های کمینه‌سازی، و مثلث‌های رو به بالا مربوط به گره‌های بیشینه‌سازی را در نظر بگیرید:



بیایید نحوه استخراج این درخت را با الگوریتم مینیماکس بررسی کنیم، با تکرار در گره‌ها با مقادیر ۳، ۱۲، و ۸، و اختصاص مقدار $3 = 12, \min(3) = 8$ به سمت چپ‌ترین کمینه‌سازی شروع می‌کنیم. سپس، $\min(2, 4, 6) = 2$ (۶ را به کمینه‌ساز میانی، و $5, 2 = \min(2)$ را به سمت راست‌ترین کمینه‌سازی، قبل از اینکه در نهایت $\max(3, 2, 2) = 3$ را به بیشینه‌ساز واقع در ریشه اختصاص دهیم، اختصاص می‌دهیم. با این حال، اگر به این وضعیت فکر کنیم، می‌توانیم به این موضوع پی ببریم که به محض بازدید از فرزند کمینه‌ساز وسط با مقدار ۲، دیگر نیازی نیست به سایر فرزندان کمینه‌ساز وسط نگاه کنیم. چراکه ما فرزندی از کمینه‌ساز وسط با مقدار ۲ دیده‌ایم، می‌دانیم که مهم نیست فرزندان دیگر چه مقداری دارند، چون مقدار کمینه‌ساز وسط حداکثر می‌تواند ۲ باشد. اکنون که این مورد مشخص شد، بیایید یک گام بعدی را بررسی کنیم. هنوز هم بیشینه‌ساز در ریشه باید از بین مقدار ۳ از کمینه‌ساز سمت چپ و مقدار $2 \geq 2$ انتخاب کند، تضمین می‌شود که بیشینه‌ساز مقدار ۳ که توسط کمینه‌ساز چپ برگردانده شده است را بر مقدار بازگشتی کمینه‌ساز وسط، بدون توجه به مقادیر فرزندان باقی‌مانده، ترجیح می‌دهد. دقیقاً به همین دلیل است که می‌توانیم درخت جست‌وجو را هرس کنیم و هرگز فرزندان باقی‌مانده کمینه‌ساز وسط را بررسی نکنیم:



اجرای چنین هرسی می تواند زمان اجرا را به $O(b^{m/2})$ کاهش دهد و به طور موثر عمق "قابل حل" ما را دو برابر کند. در عمل، اغلب بسیار کمتر از این مقدار است، اما به طور کلی می تواند جستجو در حداقل یک یا دو سطح دیگر را امکان پذیر کند. که این نکته قابل توجهی است، چراکه بازیکنی که فکر می کند ۳ حرکت جلوتر است، بر بازیکنی که فکر می کند ۲ حرکت جلوتر است برای برد ترجیح داده می شود. این هرس دقیقاً همان کاری است که الگوریتم مینیماکس با هرس آلفا-بتا انجام می دهد و به صورت زیر پیاده سازی می شود:

α : MAX's best option on path to root
 β : MIN's best option on path to root

```
def max-value(state,  $\alpha$ ,  $\beta$ ):
    initialize  $v = -\infty$ 
    for each successor of state:
         $v = \max(v, \text{value}(\text{successor}, \alpha, \beta))$ 
        if  $v \geq \beta$  return  $v$ 
         $\alpha = \max(\alpha, v)$ 
    return  $v$ 
```

```
def min-value(state,  $\alpha$ ,  $\beta$ ):
    initialize  $v = +\infty$ 
    for each successor of state:
         $v = \min(v, \text{value}(\text{successor}, \alpha, \beta))$ 
        if  $v \leq \alpha$  return  $v$ 
         $\beta = \min(\beta, v)$ 
    return  $v$ 
```

این کد را با کد اصلی مینیماکس مقایسه کنید، و توجه داشته باشید که اکنون می توانیم زودتر بدون جست و جو در میان همه جانشین ها به نتیجه برسیم.

توابع ارزیابی:

اگرچه هرس آلفا-بتا می تواند به افزایش عمقی که می توانیم اجرا کنیم، کمک می کند، اما معمولاً به اندازه کافی برای رسیدن به انتهای درخت های جستجو برای اکثر بازی ها خوب نیست. در نتیجه، ما به توابع ارزیابی روی می آوریم، توابعی که یک حالت را می گیرند و تخمینی از مقدار کمینه واقعی آن گره را به دست می آورند. به طور معمول، در این توابع حالت هایی به عنوان حالت های "بهتر" تعبیر می شوند که مقادیر بالاتری را در مقابل حالت های "بدتر" نسبت می دهند. توابع ارزیابی به طور گسترده در مینیماکس های محدود به عمق استفاده می شوند، جایی که گره های غیر پایانی واقع در حداکثر عمق قابل حل خود را به عنوان گره های پایانه در نظر می گیریم، و به آنها ابزار پایانه ساختگی می دهیم که توسط یک تابع ارزیابی با دقت انتخاب شده تعیین می شود. از آنجایی که توابع ارزیابی فقط می توانند تخمین هایی از مقادیر ابزارهای غیر پایانی ارائه دهند، تضمین بهینه بودن بازی را هنگام اجرای مینیماکس حذف می کنند.

هنگام طراحی عاملی که مینیماکس را اجرا می کند، معمولاً برای انتخاب یک تابع ارزیابی، فکر و آزمایش های زیادی انجام می شود و هر چه عملکرد تابع ارزیابی بهتر باشد، عامل به رفتار بهینه نزدیکتر می شود. علاوه بر این، قبل از استفاده از یک تابع ارزیابی، عمیق تر شدن در درخت نیز نتایج بهتری به ما می دهد، زیرا انجام محاسبات آنها در عمق درخت بازی، به خطر افتادن بهینه گی را کاهش می دهد. هدف این توابع در بازی ها بسیار مشابه توابع اکتشافی در مسائل جستجوی استاندارد است.

رایج ترین طراحی برای یک تابع ارزیابی، ترکیب خطی از ویژگی ها است.

$$Eval(s) = w_1 f_1(s) + w_2 f_2(s) + \dots + w_n f_n(s)$$

هر $f_i(s)$ مربوط به یک ویژگی استخراج شده از حالت ورودی s است و به هر ویژگی یک وزن w_i اختصاص داده می شود. ویژگی ها به سادگی برخی از عناصر یک حالت بازی هستند که می توانیم آن ها را استخراج و یک مقدار عددی به آنها اختصاص دهیم. به عنوان مثال، در یک بازی چکرز ممکن است یک تابع ارزیابی با ۴ ویژگی بسازیم: تعداد مهره های ^{۱۳}عامل، تعداد شاه های عامل، تعداد مهره های حریف و تعداد شاه های حریف. سپس وزن های مناسب را بر اساس اهمیت آن ها انتخاب می کنیم. در مثال چکرز، انتخاب وزن های مثبت برای مهره ها/پادشاهان عامل خود و وزن های منفی برای مهره ها/پادشاهان حریف بسیار منطقی است. علاوه بر این، ممکن است تصمیم بگیریم که به علت ارزش بیشتر پادشاهان نسبت به بقیه مهره ها، به ویژگی های مربوط به پادشاهان عامل/رقیب وزن های بیشتری نسبت به ویژگی های مربوط به بقیه مهره ها بدهیم. در زیر یک تابع ارزیابی احتمالی داریم که با ویژگی ها و وزن هایی که اخیراً بررسی کردیم مطابقت دارد:

$$Eval(s) = 2 \cdot agent_kings(s) + agent_pawns(s) - 2 \cdot opponent_kings(s) - opponent_pawns(s)$$

همانطور که مشخص است، طراحی تابع ارزیابی می تواند کاملاً آزاد باشد و لزوماً نباید از توابع خطی باشد. برای مثال، توابع ارزیابی غیرخطی مبتنی بر شبکه های عصبی در برنامه های یادگیری تقویتی^{۱۴} بسیار رایج هستند. مهم ترین چیزی که باید در نظر داشت این است که تابع ارزیابی تا جایی که ممکن است، باید امتیازات بالاتری را برای موقعیت های بهتر به دست بیاورد. این ممکن است نیاز به تنظیم دقیق و آزمایش های متعدد روی عملکرد عوامل با استفاده از توابع ارزیابی با ویژگی ها و وزن های مختلف داشته باشد.

¹ pawns
¹ Reinforcement Learning applications