

## مترجم: سهراب جلوه گر

ویرایش دوم، بهار ۱۳۸۸



## هوش مصنوعی

محیط  $S=20$

در  $t = 4$ :

1. 01## : 0000  $S=220$  0001
2. 00#0 : 1100  $S=208$
3. 11## : 1000  $S=196$
4. ##00 : 0001  $S=156$

محیط  $S=20$

و در  $t = 5$  - درخواست به سیستم می آید و به آخرین طبقه بندی کننده ی فعال نسبت داده می شود .

- 1) 01## : 0000  $S=220$
- 2) 00#0 : 1100  $S=208$
- 3) 11## : 1000  $S=196$
- 4) ##00 : 0001  $S=206$

محیط  $S=20$

**تولید قوانین جدید** - دسته سطل ، یک روش انتخاب قوانین و انتساب اعتبار را ارائه می نماید .

چگونه قوانین جدید را به دست بیاوریم ؟ الگوریتم ژنتیکی پایه ی ما از یک مدل جمعیت که با هم اشتراک ندارند استفاده می نماید ؛ همه ی جمعیت در زمان  $t$  ، در زمان  $t+1$  جایگزین می شود . به خوبی برای بهینه سازی عمل می کند ، ولی برای یادگیری ، زیاد مناسب نمی باشد . در عوض ، ما از نخبه سالاری<sup>۱</sup> برای حفظ

---

<sup>۱</sup> elitism

مترجم: سهراب جلوه گر

ویرایش دوم، بهار ۱۳۸۸



## هوش مصنوعی

برخی از قانون ها استفاده می نمایم ، از روش انتخاب رولت<sup>۱</sup> برای حفظ قانون ها استفاده نمایم ، البته این روش ، کندتر استنتاج می کند .

**برنامه نویسی ژنتیکی - فراتر از قوانین** - روش های تکاملی برای تولید برنامه ها به کار گرفته

شده اند . ایده ی این روش این است که یک برنامه می تواند به صورت یک S- عبارتی بیان شود ، مثلاً :

$$(x * 4 + 3)$$

ما می توانیم این را به صورت یک درخت بکشیم . ما سپس تکامل را روی یک جمعیت از برنامه ها انجام می دهیم . شایستگی ، کارایی یک برنامه برای یک عملکرد ارایه شده است . Crossover ، زیر درخت ها را با هم عوض می نماید . جهش ، عملگرها را عوض می نماید و برای تکامل استفاده می شوند ، مثال ها برنامه های فوتبال روبات ، مدارهای کنترل آنالوگ ، فیلترها و مدارهای پیچیده می باشند . در این مورد چالش ها عبارتند از ، فضای پهناور (نامحدود) برنامه ها و کمبود ساختارهای سطح بالاتر برنامه .

---

<sup>۱</sup> roulette selection



Decision trees<sup>1</sup>

مترجم: سهراب جلوه گر

ویرایش دوم، بهار ۱۳۸۸



هوش مصنوعی

در گذشته ، ما فرض کردیم که دانش اولیه توسط خبره ها به ما داده شده است و بر چگونگی استفاده از این دانش تمرکز کردیم . حال می خواهیم در این مورد صحبت کنیم که چگونه دانش را از راه مشاهده به دست آوریم و بر قانون های گزاره ای تمرکز می کنیم . مثلاً اگر هوا آفتابی و گرم باشد ، آن گاه تنیس بازی کنید و می نویسیم :

$\text{sunny} \wedge \text{warm} \rightarrow \text{PlayTennis}$

یا اگر هوا خنک باشد و بارانی باشد یا باد شدید بوزد ، آن گاه تنیس بازی نکنید و می نویسیم :

$\text{cool} \wedge (\text{rain} \vee \text{strongWind}) \rightarrow \neg \text{PlayTennis}$

## یادگیری

برای یک عامل ، یادگیری به چه معنی می باشد ؟ عامل ، دانش جدید را دریافت می نماید ، دانش جدید را به کار می گیرد ، رفتارش را تغییر می دهد و در یک کار معین ، معیار کارایی خود را بهبود می بخشد .

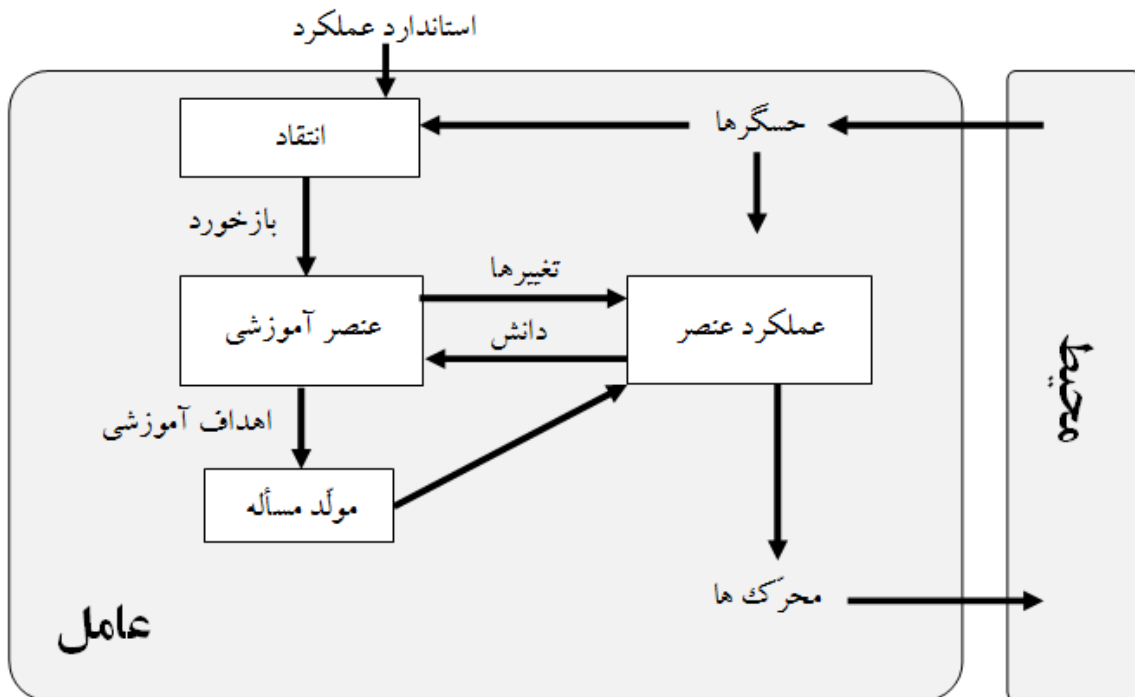
## عامل های آموزشی

مترجم: سهراب جلوه گر  
ویرایش دوم، بهار ۱۳۸۸



## هوش مصنوعی

به یاد بیاورید که ما قبلاً در مورد عامل های آموزشی صحبت کردیم :



یک عامل آموزشی دارای یک عنصر کارایی<sup>۱</sup> و یک عنصر آموزشی<sup>۲</sup> می باشد. عنصر کارایی، چیزی است که یک عامل برای تصمیم گیری در مورد این که چه کاری انجام دهد از آن استفاده می کند و این چیزی است که ما تا کنون مطالعه کرده ایم. عنصر آموزشی، چیزی است که به عامل برای بازنگری عنصر کارایی اجازه می دهد، این ممکن است به معنی اضافه نمودن یا تغییر قانون ها یا واقعیات، بازنگری

<sup>۱</sup> performance element  
<sup>۲</sup> learning element



یک مکاشفه و تغییر یک تابع جانشین<sup>۱</sup> باشد. یک عامل، برای بازنگری کردن رفتارش، به اطلاعاتی که به عامل بگوید چگونه به خوبی کارش را انجام می دهد نیازمند می باشد، این اطلاعات، بازخورد<sup>۲</sup> نام دارد.

## انواع بازخورد

در اصل سه نوع عملکرد آموزشی وجود دارد که هر کدام بازخورد متفاوتی را دارد:

**یادگیری نظارت شده<sup>۳</sup> (کنترل شده)** که در این مورد، یک منبع خارجی (که اغلب آموزش دهنده<sup>۴</sup> نام دارد)، عامل را بانمونه های برجسته زده شده<sup>۵</sup> ارایه می نماید. عامل، موارد / عملکردهای معینی را در طول طبقه بندی اشان می بیند.

**یادگیری بدون نظارت<sup>۶</sup> (کنترل نشده)** که در این مورد، آموزش دهنده ای برای ارائه ی نمونه ها وجود ندارد و عامل، معمولاً برای پیدا کردن الگوهایی در داده ها تلاش می نماید.

**یادگیری تقویتی<sup>۷</sup>** که این، یک نوع مخصوص از آموزش است که در آن عامل فقط یک پاداش را برای انجام یک عمل دریافت می نماید. ممکن است نداند یک پاداش بهینه چگونه می باشد و "بهترین" عملکرد را برای انجام نمی داند.

## یادگیری نظارت شده (کنترل شده)

<sup>۱</sup> Successor function

<sup>۲</sup> feedback

<sup>۳</sup> Supervised learning

<sup>۴</sup> teacher

<sup>۵</sup> labeled examples

<sup>۶</sup> unsupervised learning

<sup>۷</sup> reinforcement learning



یکی از عمومی ترین شکل های آموزش می باشد . عامل با مجموعه ای از داده های برچسب زده شده ارایه می شود که باید از این داده ها برای تشخیص قانون های کلی تر استفاده نماید . نمونه ها ، لیست بیماران و ویژگی ها ، چه عامل هایی مرتبط با سرطان می باشند ؟ چه عواملی ، فردی را دارای خطر می داند ؟ بهترین سؤالات برای طبقه بندی حیوانات چیست ؟ صورت چه کسی در این تصویر می باشد ؟ این پردازش یادگیری قانون های کلی از واقعیات مشخص ، استنتاج<sup>۱</sup> نام دارد .

**تعریف یک مسأله ی یادگیری** - ما می توانیم مسأله ی یادگیری را با تخمین ، به صورت یک تابع  $f$  که به ما می گوید چگونه یک مجموعه از ورودی ها را طبقه بندی نماییم تفسیر کنیم . یک مثال در این مورد یک مجموعه از ورودی های  $x$  و  $f(x)$  متناظر می باشد :

<<Mammal,Eats-Meat,Black-Stripes,Tawny>,Tiger>

ما می توانیم یک عملکرد آموزشی را به این صورت ، تعریف نماییم : یک مجموعه ی ارایه شده از نمونه های  $f$  ، یک تابع  $H$  که  $f$  را برای مثال ما تخمین می زند پیدا نمایید ،  $H$  ، فرض<sup>۲</sup> نام دارد .

**استنتاج** - ما ترجیح می دهیم  $H$  را تعمیم بدهیم ، این به این معنی است که  $H$  به درستی نمونه های دیده نشده را طبقه بندی می کند . در صورتی که فرض به درستی بتواند همه ی نمونه های آموزشی را طبقه بندی نماید ، ما آن را فرض سازگار<sup>۳</sup> می نامیم . هدف ، پیدا کردن یک فرض سازگار که در نمونه های دیده نشده هم به خوبی کار کند . ما می توانیم یادگیری را به صورت جستجو در میان یک فضا از فرض ها تصور نماییم .

<sup>۱</sup> induction

<sup>۲</sup> hypothesis

<sup>۳</sup> consistent hypothesis



**بایاس استنتاجی** - توجه نمایید که استنتاج، بی عیب نمی باشد. در انتخاب یک فرض، ما یک گمان آموزش داده شده را به وجود می آوریم. روشی که ما با استفاده از آن، این گمان را به وجود می آوریم بایاس نام دارد. نمونه ها، اصل اکام<sup>۱</sup>، معین ترین فرض، کلی ترین فرض و تابع خطی می باشند.

**مشاهده ی داده ها** - عامل ها، ممکن است دارای ابزار مختلف مشاهده کننده ی نمونه های یک فرض باشند. یک الگوریتم یادگیری دسته ای<sup>۲</sup>، یک مجموعه ی بزرگ داده ها که همه با هم ارایه می شوند و یک فرضیه ی منفرد را انتخاب می کنند، می باشد. یک الگوریتم یادگیری افزایشی، نمونه های همه در یک زمان را دریافت می کند و دائماً فرض آن را بازنگری می نماید؛ دسته معمولاً با دقت تر می باشد، اما افزایشی ممکن است با محیط عامل، مناسب تر باشد. یک عامل یادگیری فعال می تواند نمونه ها را انتخاب نماید. یک عامل یادگیری غیرفعال<sup>۳</sup>، دارای نمونه های ارایه شده با آن به وسیله ی یک منبع خارجی هستند می باشد؛ یادگیری فعال توانمندتر می باشد، اما ممکن است با محدودیت های دامنه، مناسب نباشد.

**درخت های تصمیم گیری یادگیری** - درخت های تصمیم گیری، ساختمان داده هایی هستند که یک عامل را با ابزارهای طبقه بندی کننده ی نمونه ها ارایه می نمایند. در هر گره در درخت، یک ویژگی، تست می شود.

---

<sup>۱</sup> Occam's razor، برای دو توصیف ارایه شده، آسان ترین آن ها را انتخاب نمایید - توضیحی که به مفروضات کم تری

نیاز دارد (لغت نامه ی وب برنامه ی babylon)

<sup>۲</sup> batch

<sup>۳</sup> passive

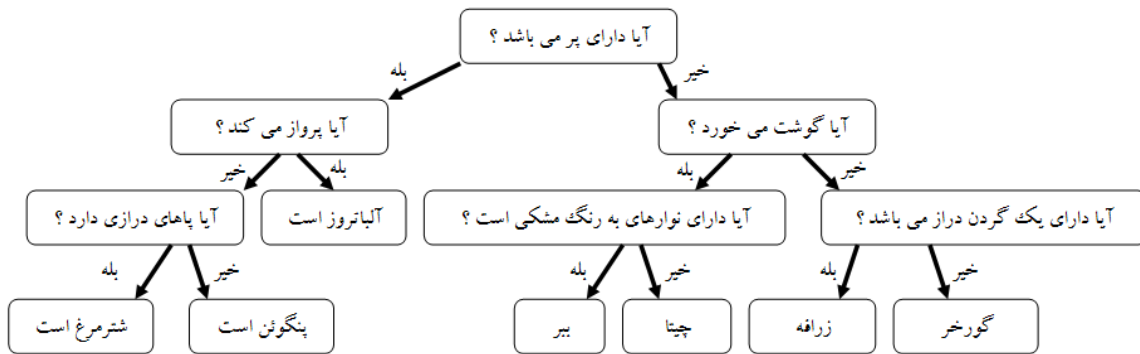


## مترجم: سهراب جلوه گر

ویرایش دوم، بهار ۱۳۸۸



## هوش مصنوعی



## اطلاعات

در مورد دنیای جاروبرقی، اتاق ها ممکن است تمیز یا کثیف باشند که برای بیان این مورد به یک بیت نیاز داریم. اگر یک اتاق دارای چهار وضعیت باشد یا هشت وضعیت چه طور؟ اگر یک اتاق فقط دارای یک وضعیت باشد چه طور؟ به چند بیت برای بیان نیاز داریم؟

## تئوری اطلاعات

به صورت رسمی تر، بیابید بگوییم که  $n$  پاسخ ممکن  $v_1, v_2, \dots, v_n$  برای یک سؤال وجود دارد و هر جواب دارای احتمال رویداد  $P(v_n)$  می باشد. محتوای اطلاعات<sup>۱</sup> جواب برای سؤال به این صورت می باشد:

$$I = \sum_{i=1}^n -P(v_i) \log_2 P(v_i)$$

برای یک سگه، این برابر است با:

$$-\frac{1}{2} \log_2 \frac{1}{2} + -\frac{1}{2} \log_2 \frac{1}{2} = 1$$

<sup>۱</sup> information content



سؤالاتی با یک احتمال زیاد جواب، دارای محتوای اطلاعاتی کمی خواهند بود. (در صورتی که سگه در یک بار ۹۹/۱۰۰ شیر بیاید،  $I=0.08$  خواهد بود.)، محتوای اطلاعاتی، برخی اوقات بی نظمی (آنروپی<sup>۱</sup>) نامیده می شود و اغلب در فشرده سازی و الگوریتم های انتقال داده ها استفاده می شود.

### استفاده از تئوری اطلاعات

برای درخت های تصمیم گیری، ما می خواهیم بدانیم ارزشمندی هر مورد ممکن چگونه می باشد، یا چه میزان از اطلاعات را نتیجه می دهد. ما می خواهیم احتمال جواب های ممکن از یک مجموعه ی مورد آموزش را تخمین بزنیم. معمولاً، یک مورد تکی برای به طور کامل مجزا نمودن نمونه های مثبت و منفی کافی نمی باشد. در عوض، ما نیاز داریم در مورد این فکر کنیم که چه مقدار بهتری بعد از سؤال وجود خواهد داشت و این، سود اطلاعات<sup>۲</sup> نام دارد.

### سود اطلاعات

در صورتی که یک مجموعه ی مورد آزمایش دارای  $p$  نمونه ی مثبت و  $n$  نمونه ی منفی باشد، بی نظمی یا آنروپی به صورت زیر خواهد بود:

$$I\left(\frac{p}{p+n}, \frac{n}{p+n}\right) = -\frac{p}{p+n} \log_2 \frac{p}{p+n} - \frac{n}{p+n} \log_2 \frac{n}{p+n}$$

ما می خواهیم خصوصیتی که برای مجزاً کردن نمونه های مثبت و منفی، نزدیک ترین می باشد را پیدا نماییم. ما با محاسبه ی باقی مانده شروع می کنیم - این اطلاعاتی است که هنوز در داده ها بعد از این که ما ویژگی  $A$  را تست نمودیم وجود دارد. می گوییم ویژگی  $A$  می تواند مقادیر ممکن  $v$  را به کار بگیرد. این تست،

<sup>۱</sup> entropy

<sup>۲</sup> information gain

## مترجم: سهراب جلوه گر

ویرایش دوم، بهار ۱۳۸۸



## هوش مصنوعی

$V$  زیر مجموعه ی جدید از داده ها را به نام  $E_1, E_2, \dots, E_v$  خواهد ساخت. باقی مانده، مجموع اطلاعات موجود در هر کدام از این زیرمجموعه ها می باشد.

$$\text{Re mainder } (A) = \sum_{i=1}^v \frac{p_i + n_i}{p + n} * I\left(\frac{p_i}{p_i + n_i}, \frac{n_i}{p_i + n_i}\right)$$

سود اطلاعات سپس می تواند به صورت تفاوت میان اطلاعات اصلی (قبل از آزمایش) و اطلاعات جدید (بعد از آزمایش) بیان شود.

$$\text{Gain}(A) = I\left(\frac{p}{p + n}, \frac{n}{p + n}\right) - \text{Re mainder } (A)$$

**مکاشفه:** ویژگی با بیش ترین سود اطلاعات را انتخاب نمایید.

**سؤال:** این کدام نوع از جستجو می باشد؟

**مثال:**

| روز | پیش بینی | دما   | رطوبت  | باد   | تنیس بازی |
|-----|----------|-------|--------|-------|-----------|
| D1  | آفتابی   | گرم   | بالا   | ملايم | نه        |
| D2  | آفتابی   | گرم   | بالا   | شدید  | نه        |
| D3  | ابری     | گرم   | بالا   | ملايم | بله       |
| D4  | بارانی   | متوسط | بالا   | ملايم | بله       |
| D5  | بارانی   | خنک   | معمولی | ملايم | بله       |
| D6  | بارانی   | خنک   | معمولی | شدید  | نه        |
| D7  | ابری     | خنک   | معمولی | شدید  | بله       |
| D8  | آفتابی   | متوسط | بالا   | ملايم | نه        |

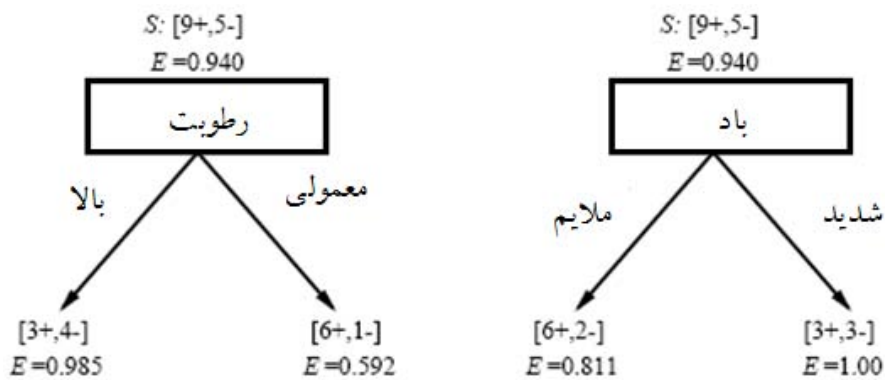
مترجم: سهراب جلوه گر  
ویرایش دوم، بهار ۱۳۸۸



## هوش مصنوعی

|     |        |       |        |       |     |
|-----|--------|-------|--------|-------|-----|
| D9  | آفتابی | خنک   | معمولی | ملايم | بله |
| D10 | بارانی | متوسط | معمولی | ملايم | بله |
| D11 | آفتابی | متوسط | معمولی | شدید  | بله |
| D12 | ابری   | متوسط | بالا   | شدید  | بله |
| D13 | ابری   | گرم   | معمولی | ملايم | بله |
| D14 | بارانی | متوسط | بالا   | شدید  | نه  |

کدام ویژگی، بهترین طبقه بندی می باشد؟



$$\text{Gain}(S, \text{رطوبت}) = 0.940 - (7/14) \cdot 0.985 - (7/14) \cdot 0.592 = .151$$

$$\text{Gain}(S, \text{باد}) = .940 - (8/14) \cdot 0.811 - (6/14) \cdot 1.0 = .048$$

## اغتشاش یا پارازیت<sup>۱</sup>

Noise<sup>۱</sup>

**مترجم: سهراب جلوه گر**

**ویرایش دوّم، بهار ۱۳۸۸**



**هوش مصنوعی**

در صورتی که دو نمونه ی دارای ویژگی های یکسان ولی مقدارهای مختلف وجود داشته باشد ، یک درخت تصمیم گیری قادر به طبقه بندی آن ها به صورت مجزاً نخواهد بود و در این صورت می گوییم که این داده ها دارای اغتشاش می باشند .

**ورودی های پیوسته ی مقداردهی شده - درخت های تصمیم گیری همچنین می توانند**

برای کار با ورودی های صحیح و مقداردهی شده به صورت پیوسته توسعه داده شوند . در این مورد برای پیدا کردن مقداری که بالاترین سود اطلاعات را نتیجه می دهد از جستجوی تپه نوردی استفاده نماید .

مترجم: سهراب جلوه گر  
ویرایش دوم، بهار ۱۳۸۸



هوش مصنوعی

## فصل بیست و چهارم

# یادگیری Q ای<sup>۱</sup>

Q - learning<sup>۱</sup>



ما به محیط های تصادفی جذب شده ایم. ما می توانیم این مسأله را به صورت یک پروسه ی تصمیم گیری مارکوفی فرمول بندی نماییم؛ مسأله دارای یک حالت اولیّه ی  $s_0$ ، یک مجموعه ی گسسته از وضعیّت ها و عملکردها، یک مدل انتقالی  $T(s,a,s')$  که احتمال رسیدن وضعیّت  $s'$  از  $s$  در موقع انجام عمل  $a$  می باشد و یک تابع پاداش  $(R(s))$  می باشد.

## روش ها<sup>۱</sup>

یک روش، راه حلی برای یک پردازش تصمیم گیری مارکوفی می باشد و عملکرد بهینه را برای انجام در هر وضعیّت، مشخص می نماید. عملکرد بهینه، عملکردی است که سودمندی مورد انتظار کاهش داده شده از این وضعیّت را بیشینه می نماید. بنابراین، سودمندی یک وضعیّت، پاداشی برای آن وضعیّت، به اضافه ی پاداش کم شده ی مورد انتظار برای دنبال کردن روش بهینه از این وضعیّت می باشد. داریم:

$$U(s) = R(s) + \gamma \max_a \sum_{s'} T(s,a,s')U(s')$$

polices<sup>۱</sup>



که عبارت فوق معادله ی بلمن<sup>۱</sup> نام دارد .

## تکرار مقدار<sup>۲</sup>

به صورت مستقیم نمی توان معادله ی بلمن را حل نمود . ما می توانیم از تکرار مقدار برای حل یک سیستم معادلات بلمن استفاده نماییم : به سودمندی های غیر پایانی ، مقدارهای تصادفی نسبت دهید . برای هر وضعیت ، عملکرد بهینه ی ارایه کننده ی این مقادیر را محاسبه نمایید . از این روش برای به روز رسانی سودمندی تخمین زده شده ی وضعیت استفاده نمایید و این کار را تا زمان نزدیک شدن به وضعیت مطلوب (همگرایی) ، ادامه دهید .

## تکرار روش<sup>۳</sup>

در روش تکرار مقدار ؛ حتی در زمانی که روش تغییر نمی کند ، همگرایی ، کند می باشد و ما واقعاً نمی توانیم نسبت به سودمندی یک وضعیت ، بی تفاوت باشیم . تکرار روش ، روش ها را به صورت مستقیم به روز رسانی می نماید .

## یادگیری یک روش

تا حالا ، ما فرض کردیم که دارای دانش زیادی هستیم . در حالت به خصوص ، ما فرض کرده ایم که مدلی از جهان شناخته شده است ؛ این ، مدل انتقال وضعیت می باشد . اگر ما یک مدل نداشته باشیم چطور ؟ همه ی چیزی که ما می دانیم این است که مجموعه ای از وضعیت ها و یک مجموعه از عملکردها وجود دارند . ما هنوز می خواهیم یک روش بهینه را یاد بگیریم .

---

<sup>۱</sup> Bellman equation

<sup>۲</sup> Value Iteration

<sup>۳</sup> Policy Iteration





## یادگیری Q ای

آموزش یک روش به صورت مستقیم، مشکل می باشد. مسأله این است که داده های ما به یک شکل نمی باشد: <وضعیت، عملکرد>. در عوض، به شکل  $R, s_1, s_2, s_3, \dots$  می باشد. از آنجایی که ما تابع انتقال را نمی دانیم، مشکل است که سودمندی یک وضعیت را یاد بگیریم. در عوض، ما یک تابع  $Q(s, a)$  را یاد خواهیم گرفت که "سودمندی" انجام عملکرد  $a$  را در وضعیت  $s$ ، تخمین می زند. به طور صریح تر،  $Q(s, a)$  مقدار  $a$  را در وضعیت  $s$  را ارایه می کند و بعد از این به صورت بهینه عمل می کند.

$$Q(s, a) = R(s, a) + \gamma \max_{a'} \sum_{s'} T(s, a, s') U(s')$$

روش بهینه این است که عملکرد با بالاترین مقدار  $Q$  را در هر وضعیت، به دست آوریم. اگر عامل بتواند  $Q(s, a)$  را یاد بگیرد، می تواند عملکردهای بهینه را حتی بدون دانستن تابع پاداش یا تابع انتقال داشته باشد.

## یادگیری تابع Q

برای یادگیری  $Q$ ، ما باید بتوانیم مقدار یک عملکرد را در یک وضعیت، حتی اگر پاداش ها در طول زمان پخش شده اند تخمین بزنیم. ما می توانیم این کار را به صورت تکراری انجام دهیم. توجه کنید که  $U(s) = \max_a Q(s, a)$ . سپس ما می توانیم معادله را برای  $Q$  به صورت زیر بازنویسی نماییم:

$$Q(s, a) = R(s, a) + \gamma \max_{a'} Q(s', a')$$

بیاید تخمین  $Q(s, a)$  را به صورت  $\hat{Q}(s, a)$ ، مشخص نماییم. ما یک جدول لیست کننده ی هر جفت وضعیت - عملکرد و مقدار تخمین زده شده ی  $Q$  را نگهداری می نماییم. عامل، وضعیت  $s$  را درک می کند، عملکرد  $a$  را انتخاب می کند و سپس درک می کند که پاداش برابر با  $r = R(s, a)$  می باشد، که



آن را دریافت می کند و وضعیت جدید  $s'$  می باشد. سپس جدول  $Q$  را با توجه به فرمول زیر به روز می کند :

$$\hat{Q}(s, a) = r + \gamma \max_{a'} \hat{Q}(s', a')$$

عامل، تخمینی از  $\hat{Q}$  را برای  $s'$  برای تخمین  $\hat{Q}$  برای  $s$ ، استفاده می نماید. توجه کنید که عامل نیاز به هیچ دانشی در مورد  $R$  یا تابع انتقال برای اجرای این کار ندارد. آموزش  $Q$  ای برای نزدیک شدن (همگرایی) ضمانت شده است به شرطی که پاداش ها محدود شده باشند، عامل جفت وضعیت - عملکرد های این شکلی را به روشی که معمولاً نامحدود می باشد انتخاب می کند. این به این معنی است که یک عامل باید دارای احتمالی غیر صفر از انتخاب هر  $a$  در هر  $s$  به صورت رشته ای از جفت های وضعیت - عملکرد با نگرش نامحدود باشد.

## اکتشاف

بنابراین چگونه آن را به وجود بیاوریم؟ یادگیری  $Q$  ای دارای یک تفاوت مشخص از دیگر الگوریتم هایی که تا کنون دیده ایم می باشد و آن این است که عامل می تواند عملکردها و دیدی که آن ها به دست می آورند را انتخاب نماید و این یادگیری پویا نام دارد. عامل، مایل به بیشینه کردن عملکرد می باشد؛ این به این معنی است که کاری را که در حال حاضر به نظر می رسد بهترین باشد را انجام می دهد، اما اگر عامل هیچ گاه عملکردهای با "دید بد" نداشته باشد، نمی تواند از خطاها بیرون بیاید.  $\hat{Q}$  خیلی بادقت نمی باشد، بنابراین، عملکردهای غیر بهینه را امتحان خواهیم کرد. از این پس، در صورتی که  $\hat{Q}$  بهتر شود، ما عملکردهای بهینه را انتخاب خواهیم کرد.

active learning<sup>1</sup>



## اکتشاف بولتزمان<sup>۱</sup>

یک راه برای انجام این کار استفاده از اکتشاف بولتزمان است. ما یک عمل با احتمال زیر را انجام

می دهیم:

$$P(a | s) = \frac{k^{\hat{Q}(s,a)}}{\sum_j k^{\hat{Q}(s,a_j)}}$$

که  $k$  پارامتر دما می باشد. این فرمول، شبیه به فرمولی است که ما در روش شبیه سازی گرم و سرد کردن استفاده کردیم.

## یادگیری مستحکم<sup>۲</sup>

یادگیری  $Q$  ای نمونه ای از آن چه که ما یادگیری مستحکم می نامیم می باشد. در یادگیری مستحکم، عامل ها نمونه هایی از این که چگونه عمل نمایند را نمی بینند. در عوض، آن ها عملکردها را انتخاب می نمایند و پاداش ها یا مجازات ها را دریافت می نمایند. یک عامل می تواند یاد بگیرد، حتی در زمانی که عمل غیر بهینه ای را انجام می دهد. تعمیم ها به پاداش تأخیر داده شده و پردازش های غیرقطعی تصمیم گیری مارکوفی رسیدگی می کنند.

## خلاصه

یادگیری  $Q$  ای برخی اوقات به صورت یادگیری مستقل از مدل<sup>۱</sup> شناخته می شود. در این روش، عامل فقط نیاز به انتخاب عملکردها و دریافت پاداش ها دارد. مشکلات این روش، چگونه تعمیم بدهد؟،

---

<sup>۱</sup> Boltzmann exploration  
<sup>۲</sup> reinforcement learning

مترجم: سهراب جلوه گر

ویرایش دوم، بهار ۱۳۸۸



هوش مصنوعی

مقیاس پذیری<sup>۲</sup> و سرعت نزدیک شدن (همگرایی) می باشد. یادگیری Q ای برای یک الگوریتم مفید تجربی به سطح مطلوبی رسیده است.

---

model-free learning<sup>۱</sup>  
scalability<sup>۲</sup>

مترجم: سهراب جلوه گر  
ویرایش دوم، بهار ۱۳۸۸



هوش مصنوعی

## فصل بیست و پنجم

# برنامه ریزی<sup>۱</sup>

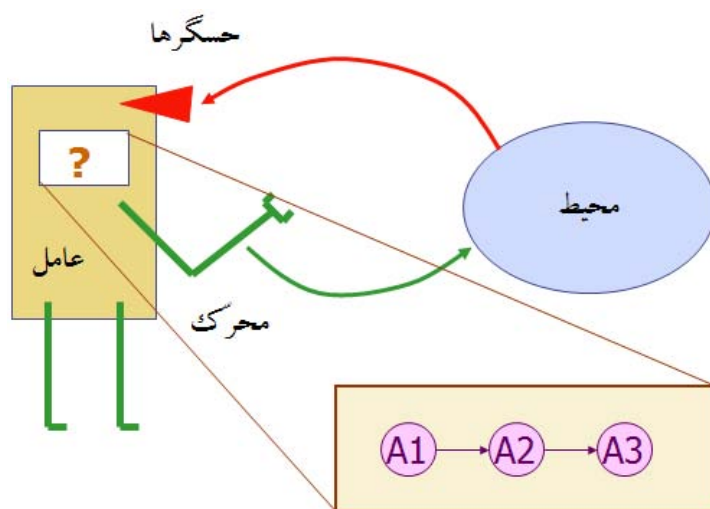
planning<sup>۱</sup>

مترجم: سهراب جلوه گر  
ویرایش دوم، بهار ۱۳۸۸



هوش مصنوعی

عامل برنامه ریزی



برنامه ریزی

مترجم: سهراب جلوه گر  
 ویرایش دوّم، بهار ۱۳۸۸

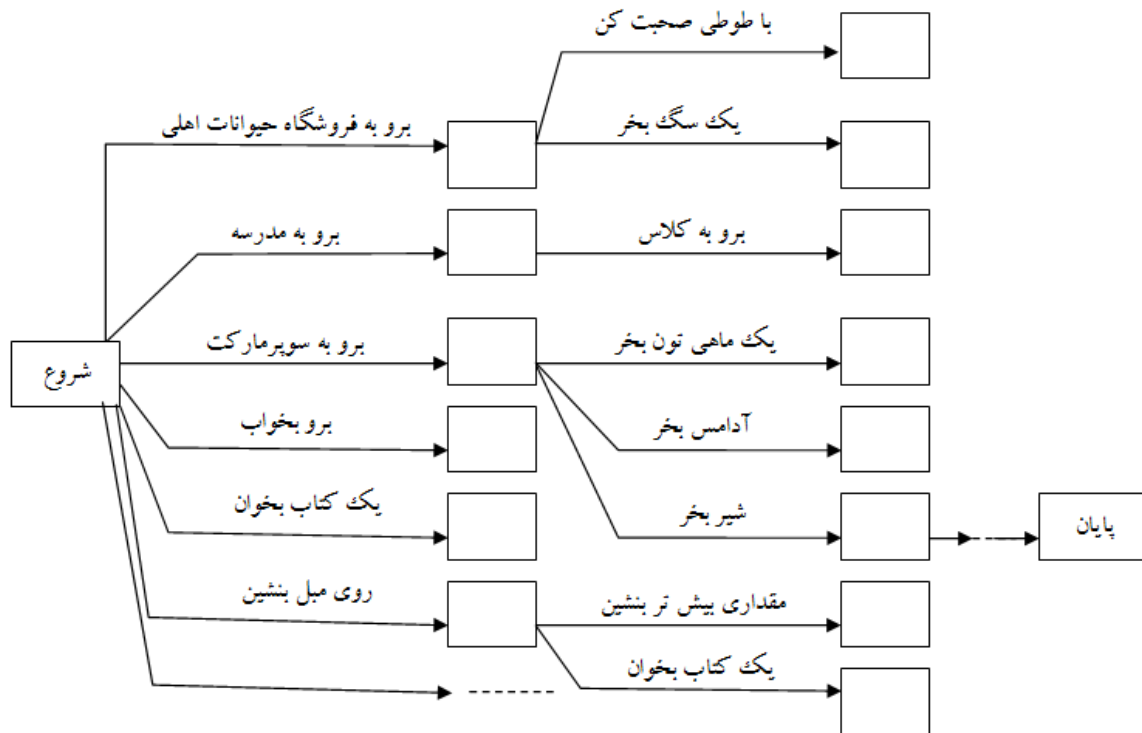


## هوش مصنوعی

برنامه ریزی، تعبیه کردن یک رشته از عملکردها برای دسترسی به یک هدف می باشد. برنامه ریزی کلاسیک، کاملاً قابل مشاهده، قطعی، محدود، ثابت<sup>۱</sup> و گسسته است. زبان مشخصی برای آرایه ی مسایل برنامه ریزی وجود دارد.

### جستجو و برنامه ریزی

الگوریتم های جستجوی استاندارد به نظر می رسد که به سختی شکست بخورند:



<sup>۱</sup> static

مترجم: سهراب جلوه گر

ویرایش دوم، بهار ۱۳۸۸



## هوش مصنوعی

مشکلاتی که در روش جستجو وجود دارد عبارتند از: تعداد زیادی عملکردهای نامربوط وجود دارد، پیدا کردن مکاشفه های خوب مشکل می باشد و از تجزیه ی مسأله نمی توان سودی برد. به صورت خودکار مکاشفه ها از ارایه مشتق می شوند؛ مانند مسأله ی ارضای محدودیت.

### مقایسه ی برنامه ریزی با منطق

| برنامه ریزی                        | جستجو                 |          |
|------------------------------------|-----------------------|----------|
| عبارت های منطقی                    | ساختمان داده های جاوا | وضعیت ها |
| پیش شرط ها <sup>۱</sup> / نتیجه ها | کد جاوا               | عملکردها |
| عبارت منطقی                        | کد جاوا               | هدف      |
| محدودیت های روی عملکردها           | رشته ای از $S_0$      | برنامه   |

### زبانی برای مسایل برنامه ریزی

حل کننده ی مسأله ی مؤسسه ی تحقیق استنفورد<sup>۲</sup> (استریپس): جهان توسط شرط های منطقی توصیف می شود، وضعیت به صورت اجتماع حرف های مثبت می باشد و می تواند به صورت

<sup>۱</sup> preconditions

<sup>۲</sup> STanford Research Institute Problem Solver (STRIPS)



مترجم: سهراب جلوه گر

ویرایش دوم، بهار ۱۳۸۸



هوش مصنوعی

گزاره ای؛ به عنوان مثال،  $Happy \wedge Hungry$  برای نمایش وضعیت عامل و یا منطق مرتبه ی اول و عبارت های مستقل از تابع باشد؛ به عنوان مثال،

$$At(Plane_1, Verona) \wedge At(Plane_2, Malpensa)$$

**فرض جهان بسته**؛ هر شرطی که نام برده نشده غلط می باشد.

**هدف**، وضعیتی است که به صورت جزئی مشخص شده است. اگر محدودیت ها همه لفظ هایی از هدف باشند یک وضعیت یک هدف را ارضا می نماید. به عنوان مثال،

$$At(Plane_1, Verona) \wedge At(Plane_2, Malpensa)$$

هدف  $At(Plane_2, Malpensa)$  را ارضا می نماید.

## عملکردهای استریپس

عملکردها توسط موارد زیر مشخص می شوند:

**پیش شرط ها**<sup>۱</sup>: زمانی را مشخص می کنند که عملکرد می تواند به کار گرفته شود.

**اثرات**: وضعیت توسط عملکردها تغییر می یابد. در این مورد، لیست اضافه، گزاره هایی هستند که صحیح از آب در می آیند و لیست حذف، گزاره هایی که غلط از آب در می آیند.

عملکردها شامل متغیرها می باشند، یک طرح عملکرد منفرد عملکردهای مختلف را نمایش می دهد (نمونه ای از متغیرها). توصیف عملکردها به صورت منظم می باشد و زبان، محدود شده است.

**طرح عملکرد:**

---

<sup>۱</sup> preconditions



$At(p) \wedge Sells(p, x)$

Buy(x)

Have(x)

Buy(x)

پیش شرط ،  $At(p) \wedge Sells(p, x)$  می باشد و اثر ،  $Have(x)$  است و باید به این نکته توجه کنیم که ، هیچ اطلاعاتی در مورد چگونگی اجرای عملکرد وجود ندارد! . چون که زبان محدود شده است پس الگوریتم کارا می باشد . عملکرد ، نام و لیست پارامتر را مشخص می نماید . پیش شرط ، اجتماع لفظ های مثبت و اثر ، اجتماع لفظ ها (مثبت یا منفی) است . یک مجموعه ی کامل از عملگرهای استریس می تواند به یک مجموعه از اصول جانشین وضعیّت<sup>۱</sup> ترجمه شود .

## معانی<sup>۲</sup>

با داشتن یک وضعیّت (اجتماع لفظ ها) ، اگر انتساب یک متغیر ، در ارتباط با لفظ های شامل شده ی وضعیّت باشد ، پیش شرط ارضا می شود ؛ به عنوان مثال ، وضعیّت  $At(HW) \wedge Sell(HW, Drill)$  ، پیش شرط  $At(p) \wedge Sell(p, x)$  را با انتساب  $p/HW$  و  $x/Drill$  ارضی می کند .

عملکردهای با پیش شرط های ارضا شده می توانند در این موارد به کار گرفته شوند : حذف عناصر از لیست حذف ، اضافه نمودن عناصر به لیست اضافه ، وضعیّت جدید :

$$At(HW) \wedge Sell(HW, Drill) \wedge Have(Drill)$$

مثال : خرید

<sup>۱</sup> successor – state axioms

<sup>۲</sup> semantics

مترجم: سهراب جلوه گر  
ویرایش دوّم، بهار ۱۳۸۸



## هوش مصنوعی

عملکردها عبارتند از :

$Buy(x)$

پیش شرط :  $At(store)$  و  $Sells(store,x)$

اثر :  $Have(x)$

$Go(x,y)$

پیش شرط :  $At(x)$

اثر :  $At(y)$  و  $\neg At(x)$

شروع :

$At(Home) \wedge Sells(SM, Milk) \wedge Sells(SM, Banana) \wedge Sells(HWS, Drill)$

هدف :

$Have(Milk) \wedge Have(Banana) \wedge Have(Drill)$

مثال : انتقال محموله ی هوایی

عملکردها :

$Load(c,p,a)$

پیش شرط :  $At(c,a) \wedge At(p,a) \wedge Cargo(c) \wedge Plane(p) \wedge Airport(a)$

اثر :  $\neg At(c,a) \wedge In(c,p)$

$Unload(c,p,a)$

مترجم: سهراب جلوه گر

ویرایش دوم، بهار ۱۳۸۸



هوش مصنوعی

پیش شرط :  $In(c, p) \wedge At(p, a) \wedge Cargo(c) \wedge Plane(p) \wedge Airport(a)$

اثر :  $At(c, a) \wedge \neg In(c, p)$

$Fly(p, from, to)$

پیش شرط :  $At(p, from) \wedge Plane(p) \wedge Airport(from) \wedge Airport(to)$

اثر :  $\neg At(p, from) \wedge At(p, to)$

عملکردها عبارتند از Load ، Unload و نتیجه ی Fly . گزاره ها ،  $In(0,0)$  و  $At(0,0)$  می باشند . انواع گزاره ها هم  $Cargo(0)$  ،  $Plane(0)$  و  $Airport(0)$  هستند .

شروع :

$At(C_1, SFO) \wedge At(C_2, JFK) \wedge At(P_1, SFO) \wedge At(P_2, JFK)$   
 $\wedge Cargo(C_1) \wedge Cargo(C_2) \wedge Plane(P_1) \wedge Plane(P_2)$   
 $\wedge Airport(JFK) \wedge Airport(SFO)$

هدف :  $At(C_1, JFK) \wedge At(C_2, SFO)$

یک راه حل به صورت زیر می باشد :

$Load(C_1, P_1, SFO), Fly(P_1, SFO, JFK), Unload(C_1, P_1, JFK),$   
 $Load(C_2, P_2, JFK), Fly(P_2, JFK, SFO), Unload(C_2, P_2, SFO)$

## برنامه ریزی استریپس

مسأله ی برنامه ریزی استریپس : یک رشته از عملکردهایی که به هدف منجر می شوند را پیدا نمایید . وضعیت ها و هدف ها توسط یک اجتماع از لفظ ها تعریف می شوند .

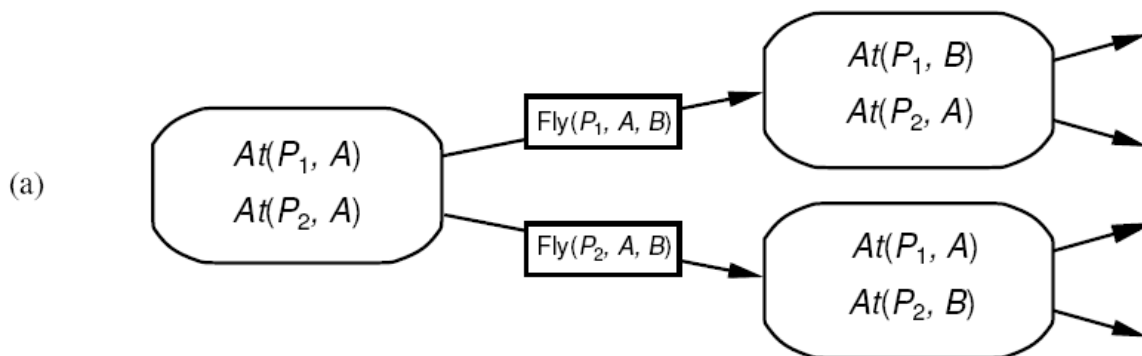


**فضای حالت جستجو:** جستجوی مستقیم از وضعیت اولیه تلاش می کند تا به هدف برسد. جستجوی معکوس از هدف تلاش و طرح ریزی می کند که به وضعیت ابتدایی برسد.

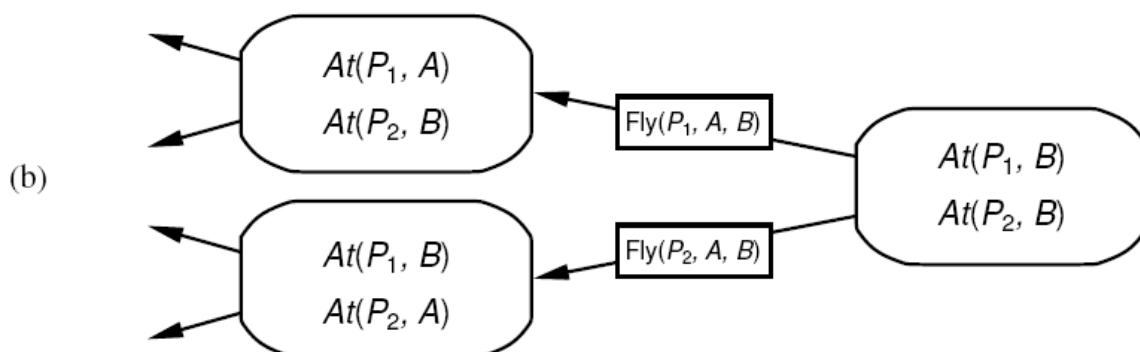
**فضای برنامه ی جستجو:** برنامه ریزی با مرتبه ی جزئی<sup>۱</sup>، فضای جزئی ساخته شده توسط طرح ها را جستجو می کند.

### فضای حالت جستجو

مسأله ی برنامه ریزی، مسأله ی جستجو را تعریف می کند. وضعیت اولیه، وضعیت شروع می باشد، آزمون هدف، بررسی می کند که آیا هدف برآورده (ارضا) می شود یا نه، عملکردها، عملگرها را تعریف می کنند، هزینه مرحله، معمولاً ۱ می باشد. در زیر، شکل a، جستجوی مستقیم و شکل b، جستجوی معکوس می باشد:



<sup>۱</sup> Partial-Order Planning(POP)



**جستجوی مستقیم** - حلقه ی جستجوی اصلی: یک عملکرد را انتخاب نمایید و پیش شرط را با وضعیت یکی نمایید. در صورتی که پیش شرط برآورده شد، عملکرد تولید کننده ی یک وضعیت جدید را به کار ببرید. بررسی کنید که آیا وضعیت جدید، هدف را برآورده می نماید یا نه. در صورتی که فضای حالت، محدود باشد، برای الگوریتم های جستجوی کامل، کامل می باشد. این روش به دلیل عملکردهای نامربوط، ناکارآمد می باشد و به مکاشفه های خوب، نیازمند می باشد.

**جستجوی معکوس** - ایده ی این روش این است که عملکردهای مربوط را فقط با شروع از هدف انتخاب نمایید. عملکرد، در زمانی که به یکی از متحدانش دست می یابد (لیست اضافه)، مربوط می باشد. در زمانی که هر لفظ مربوط را بی اثر<sup>۱</sup> نمی کند (لیست حذف)، سازگار می باشد.

یک عملکرد مربوط و سازگار را انتخاب نمایید و وضعیت جدید را تولید نمایید: لیست اضافه را پاک نمایید و پیش شرط ها را اضافه نمایید.

این روش با یک وضعیت برآورده شده توسط وضعیت ابتدایی خاتمه می یابد و روی وضعیت های مربوط و سازگار منشعب می شود.

### فضای حالت و فضای برنامه

<sup>۱</sup> undo

## مترجم: سهراب جلوه گر

ویرایش دوم، بهار ۱۳۸۸



## هوش مصنوعی

جستجوی مستقیم و معکوس، رشته های خطی عملکردها را اکتشاف می نماید که این کار لازم نمی باشد؛ به راه حل انتقال هوایی محموله توجه نمایید:

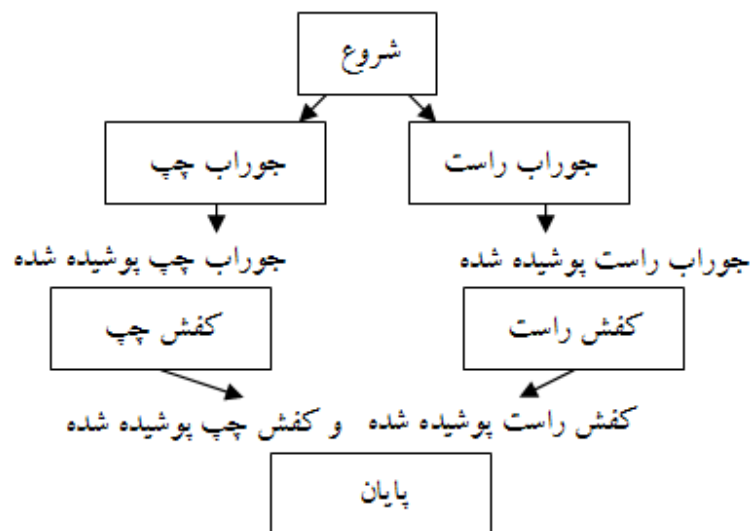
$Load(C_1, P_1, SFO), Fly(P_1, SFO, JFK), Unload(C_1, P_1, JFK),$   
 $Load(C_2, P_2, JFK), Fly(P_2, JFK, SFO), Unload(C_2, P_2, SFO)$

تنها ترتیب لازم، میان Load، Fly و Unload می باشد که توسط اثرات سببی به هم متصل شده اند؛ به عنوان مثال،  $Load(C_1, P_1, SFO)$  به  $In(C_1, P_1)$  استفاده شده توسط پیش شرط Unload دست می یابد. نیازی نیست که  $Load(C_2, P_2, JFK)$  را بعد از  $Unload(C_1, P_1, JFK)$  قرار دهیم.

**طرح با مرتبه ی جزئی<sup>۱</sup>:** گراف عملکردها شامل شروع و پایان می باشد (گراف  $\equiv$  مرتبه ی

جزئی >)

**مثال:** طرح با مرتبه ی جزئی:



<sup>۱</sup> partial-order plan



طرح های با مرتبه ی کلی :



طرح های منظم شده به صورت جزئی - مجموعه ای از مراحل با مرحله ی شروع<sup>۱</sup>

دارای توصیف وضعیّت اولیّه به صورت اثرش می باشد، مرحله ی پایانی<sup>۲</sup> دارای توصیف هدف به صورت پیش شرطش می باشد، اتصالات سببی<sup>۳</sup> از نتیجه های یک مرحله به پیش شرط دیگری به وجود می آیند و مرتب سازی زمانی<sup>۴</sup> میان جفت های مراحل می باشد.

شرط باز<sup>۵</sup>، پیش شرط یک گام که هنوز به صورت سببی متصل نشده است می باشد. طرح در صورتی کامل است که هر پیش شرط در دسترس باشد، پیش شرط در صورتی قابل دسترس می باشد که اثر یک

<sup>۱</sup> start step

<sup>۲</sup> finish step

<sup>۳</sup> causal links

<sup>۴</sup> temporal ordering

<sup>۵</sup> open condition



## هوش مصنوعی

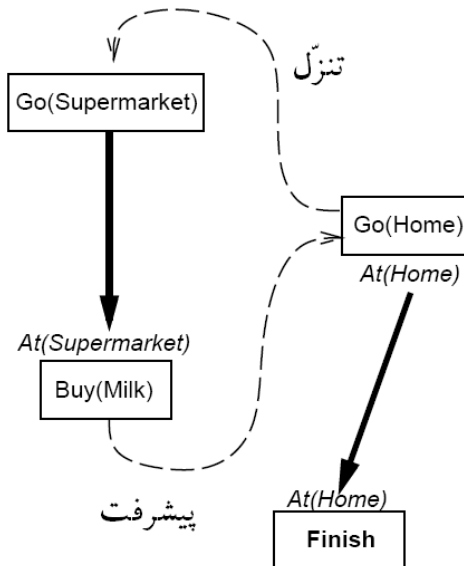


مترجم: سهراب جلوه گر

ویرایش دوم، بهار ۱۳۸۸

مرحله قبل تر باشد و هیچ مرحله ی میانی، آن را بی اثر ننماید (تناقض نداشته باشند). طرح ها نمی توانند دارای حلقه باشند.

ناسازگاری ها و پیشرفت <sup>۱</sup> / تنزل <sup>۲</sup> - یک عملکرد ناسازگار یک مرحله ی میانی بالقوه است که شرط قابل دسترس توسط اتصال سببی را خراب می نماید. به عنوان مثال، (Go(Home)، (At(Supermarket) را حذف می نماید:



تنزل: قبل از Go(Supermarket) قرار دهید

پیشرفت: بعد از Buy(Milk) قرار دهید

promotion <sup>۱</sup>

demotion <sup>۲</sup>

مترجم: سهراب جلوه گر  
ویرایش دوم، بهار ۱۳۸۸



هوش مصنوعی

## فهرست برخی از منابع و مآخذ

بیش تر مطالب استفاده شده برای ترجمه ی این کتاب ، مطالب به زبان انگلیسی موجود در پایگاه های اینترنتی دانشگاه های کشور ایالات متحده ی آمریکا می باشد . در زیر ، فهرست برخی از پایگاه های اینترنتی ای که در ترجمه ی این کتاب از آن ها استفاده شده را مشاهده می نمایید :

۱- مطالب موجود در پایگاه اینترنتی دانشگاه سانفرانسیسکو ایالات متحده ی آمریکا ، به نشانی  
[www.cs.usfca.edu](http://www.cs.usfca.edu) اینترنتی

۲- مطالب موجود در پایگاه اینترنتی دانشکده ی برکلی دانشگاه کالیفرنیا ایالات متحده ی آمریکا ،  
با آدرس اینترنتی <http://www.cs.berkeley.edu> (مطالب استوارت راسل و پیتر نورویگ)

۳- مطالب موجود در پایگاه اینترنتی دانشکده ی علوم کامپیوتر دانشگاه دولتی تگزاس ایالات متحده  
ی آمریکا ، به آدرس اینترنتی [www.cs.txstate.edu](http://www.cs.txstate.edu)



۴- مطالب موجود در پایگاه اینترنتی <http://www-formal.stanford.edu> (مطالب جان مک کارتی، استاد دانشگاه استنفورد)

۵- مطالب گردون اس. نواک جی. آر<sup>۱</sup> استاد دانشگاه تگزاس ایالات متحده ی آمریکا به نشانی اینترنتی [www.cs.utexas.edu](http://www.cs.utexas.edu)

۶- مطالب Milos Hauskrecht استاد دانشگاه شهر پیتسبورگ کشور ایالات متحده ی آمریکا به نشانی اینترنتی [www.cs.pitt.edu](http://www.cs.pitt.edu)

۷- مطالب موجود در پایگاه اینترنتی دانشگاه مرلند کشور ایالات متحده ی آمریکا با آدرس اینترنتی <http://www.cs.umd.edu>

۸- مطالب موجود در پایگاه اینترنتی مدرسه ی علوم کامپیوتر و مهندسی، به نشانی اینترنتی [www.cse.unsw.edu.au](http://www.cse.unsw.edu.au)

۹- مطالب موجود در پایگاه اینترنتی جامعه ی هوش محاسباتی، به نشانی اینترنتی

<http://ebrains.la.asu.edu/~jennie/tutorial/index.htm>

۱۰- مطالب موجود در پایگاه اینترنتی مهندسان سیاهپوست، به آدرس اینترنتی [www.nsbe.org](http://www.nsbe.org)

۱۱- مطالب موجود در پایگاه اینترنتی دانشگاه آزاد بوزن ایتالیا، به نشانی اینترنتی <http://www.inf.unibz.it>

۱۲- مطالب موجود در پایگاه اینترنتی دانشگاه کوبلنز آلمان، به نشانی اینترنتی <http://www.uni-koblenz.de>

---

<sup>۱</sup> Gordon S. Novak Jr. -

مترجم: سهراب جلوه گر  
ویرایش دوّم، بهار ۱۳۸۸



## هوش مصنوعی

۱۳- مطالب موجود در پایگاه اینترنتی دانشکده ی علوم کامپیوتر دانشگاه کورک ایرلند ، با آدرس

اینترنتی <http://www.cs.ucc.ie>

۱۴- مطالب موجود در پایگاه اینترنتی دانشکده ی علوم کامپیوتر و اطلاعات نروژ ، به آدرس

اینترنتی [www.idi.ntnu.no](http://www.idi.ntnu.no)

۱۵- مطالب موجود در پایگاه اینترنتی دانشگاه ملی تایوان ، به آدرس اینترنتی

[www.csie.ntu.edu.tw](http://www.csie.ntu.edu.tw)

۱۶- مطالب موجود در پایگاه اینترنتی دانشگاه لیورپول <sup>۱</sup> انگلستان به نشانی اینترنتی

<http://www.csc.liv.ac.uk>

۱۷- مطالب موجود در پایگاه اینترنتی دانشگاه خصوصی سینسای <sup>۲</sup> کشور ترکیه

۱۸- مطالب موجود در پایگاه اینترنتی دانشگاه شمال غربی <sup>۳</sup> کشور ایالات متّحده ی آمریکا به نشانی

اینترنتی <http://www.eecs.northwestern.edu>

۱۹- دانشکده ی علوم کامپیوتر دانشگاه بولوگنا کشور ایتالیا به نشانی اینترنتی

<http://www.cs.unibo.it>

---

<sup>۱</sup> - Liverpool

<sup>۲</sup> - Sabancı University

<sup>۳</sup> - NORTHWESTERN UNIVERSITY