

## سیستمهای خبره استقرایی و کاربرد آن در آمار

سید باقر میراشرفی\* - حسن صادقی

دانشگاه فردوسی مشهد

### خلاصه

سیستمهای خبره بر اساس هوش مصنوعی (AI) وبا استفاده از زبانهای LISP, PROLOG و با زبانهای برنامه سازی پیشرفته دیگر ساخته میشوند.

آمار و هوش مصنوعی بطور متقابل به توسعه یکدیگر کمک میکنند. از یک طرف روشهای آماری برای توصیف مفاهیم مبهم در سیستمهای خبره مورد استفاده قرار میگیرد. از طرف دیگر سیستمهای خبره بعنوان ابزار فراگیری مفاهیم میتوانند برای حل تعداد زیادی از مسائل طبقه بندی کردن آماری مورد استفاده قرار گیرد. بخصوص در روشهای تحلیل رگرسیون، تحلیل ممیزی، تحلیل خوشه ای و تحلیل سریهای زمانی کاربرد فراوانی دارد.

در این مقاله یکی از الگوریتمهایی را که کاربرد زیادی دارد مورد بررسی قرار میدهیم. این الگوریتم توسط Quinlan در سال ۱۹۷۹ تحت عنوان ID3، برای حالتی که متغیر کلاس کیفی باشد، ارائه و تا سال ۱۹۸۶ اصلاح و تکمیل گردید. نتیجه این الگوریتم قواعدی است که بر اساس آن به مسائل طبقه بندی جواب داده میشود. در کنار این الگوریتم برای حالتی که متغیرها پیوسته باشند الگوریتم Regression Tree توسط Breiman (۱۹۸۴) مطرح میشود.

## ۱- مقدمه

شامل سیستم‌های خبره استقرانی بر اساس قواعد استقراء داده‌های آماری بوده و متعلق به دسته‌ای از روش‌ها تحت عنوان "یادگیری از روی مثالها" می‌باشد. یادگیری از روی مثالها را می‌توان جزئی از تکنیک یادگیری ماشینی ( ML ) دانست. مضاف بر اینکه یادگیری ماشینی بخشی از هوش مصنوعی را تشکیل می‌دهد .

در کنار سیستم‌های خبره استقرانی، تعدادی از الگوریتم‌های شبکه عصبی و نیز روش‌های آماری کلاسیک، نظیر کرسکون، سری‌های زمانی و تحلیل سری‌وی، از طریق یادگیری از روی مثالها تحلیل می‌شوند .

سیستم‌های خبره استقرانی پیشرفت اساسی در جهت تهیه فرآیندهای اتوماتیک و فراگیری دانش‌ها داشته است. با استفاده از تجربیات گذشته که به‌صورت نمونه‌هایی در دسترس هستند، می‌توانیم قوانینی بر اساس دانش‌ها تولید کنیم. بیشتر سیستم‌های خبره استقرانی که در حال حاضر مورد استفاده قرار می‌گیرد و کاربرد تجاری نیز پیدا کرده است، بر اساس مقاله مشترک ( ۱۹۸۶ ) Quinlan و ( ۱۹۸۴ ) Breiman می‌باشد .

که در بخش بعدی به معرفی این الگوریتم‌ها و کاربرد آن می‌پردازیم .

## ۲- تحلیل سری‌وی طبقه بندی

در تحلیل سری‌وی طبقه بندی به روش کلاسیک، عموماً " بر اساس روش کمترین مربعات خطا و می‌توانیم کردن تابع زیان، و با استفاده از نمونه‌ها، به نتایج سری‌وی مشخص جهت تحلیل دسترسی پیدا می‌کنیم .

اما در روش‌های سیستم‌های خبره استقرانی بر اساس الگوریتم‌های مشخص، که توفیق یکی از این الگوریتم‌ها پیدا خواهد آمد، به یک درخت تصمیم گیری رسیده و بر

اساس درخت. قواعد تصمیم‌گیری مشخص خواهد شد. و با استفاده از قواعد بدست آمده، به تحلیل مسائل مورد نظر می‌پردازیم.

### ۳- الگوریتم ID3 (۱۹۸۶)

اولین بار در سال ۱۹۷۹ این الگوریتم توسط Quinlan مطرح که پس از اصلاح در سال ۱۹۸۶ بصورت کاملتری تحت عنوان "Intractive Dechotomizer 3" ارائه گردید. در این روش با استفاده از اطلاعات نمونه‌ها و بر اساس دستورالعمل‌های زیر به نتایج مورد نظر خواهیم رسید. فرم کلی داده‌ها در جدول شماره ۱ نشان داده می‌شود. این الگوریتم زمانی مورد استفاده قرار می‌گیرد، که مقادیر مربوط به طبقه‌ها (کلاسها) مقادیر کیفی (کست) باشند.

#### دستور العمل ID3

- ۱- ابتدا کل نمونه‌ها را بعنوان یک مجموعه S در نظر می‌گیریم، که اصطلاحاً "به آن گره نامیم".
- ۲- اگر تمام نمونه‌ها در گره S دارای کلاس مشترک C باشند، آنگاه دسته بندی مربوط به این گره C خواهد بود و به دستورالعمل بعدی می‌پردازیم.
- در غیر اینصورت خصیصه A را که دارای بیشترین اطلاعات می‌باشد انتخاب می‌کنیم. سپس S را به مجموعه‌های جدا از هم بر اساس مقادیر A افراز می‌کنیم و هر یک از این مجموعه‌های افراز شده را بعنوان گره S در نظر گرفته و دستور العمل مرحله ۲ را برای تمام این گره‌ها اجرا می‌کنیم.
- ۳- در این مرحله تمامی گره‌های نهایی کلاس یکسانی داشته و الگوریتم خاتمه پیدا می‌کند.

جدول شماره ۱ - نرم کلی داده‌ها در یادگیری از روی مثالها

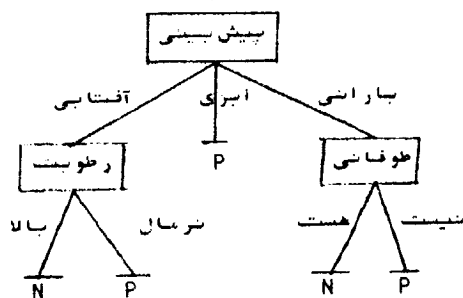
طبقه‌ها (کلاسها)	خصیصه‌ها (متغیرها)	نمونه‌ها
$C_1$	$\alpha_{11} \quad \alpha_{12} \quad \dots \quad \alpha_{1m}$	۱
$C_2$	$\alpha_{21} \quad \alpha_{22} \quad \dots \quad \alpha_{2m}$	۲
.	.	.
.	.	.
.	.	.
$C_n$	$\alpha_{n1} \quad \alpha_{n2} \quad \dots \quad \alpha_{nm}$	n

مثال ۱- در یک تحقیق وضعیت هوای صبح چهارده روزمورد مطالعه قرار گرفت. و هدف از این تحقیق این است که بر اساس مقادیر متغیرهای مورد بررسی، از مثبت بودن شرایط هوا برای یک فعالیت به نوس، مطلع شویم. متغیرها عبارتند از پیش بینی هوا با مقادیر آفتابی، ابری و بارانی و متغیر دما با مقادیر سرد، معتدل و گرم، همچنین متغیر رطوبت با مقادیر بالا و نرمال. متغیر طوفانی با مقادیر درست و غلط. اطلاعات مربوط به این تحقیق در جدول شماره ۲ آمده است.

جدول شماره ۲ - یک مجموعه اجرایی کوچک

شماره ها	خصیصه ها (متغیرها)				کلاس ها
	پیش بینی ها	دما	رطوبت	طوفانی	
۱	آفتابی	گرم	بالا	نیست	N
۲	"	"	"	هست	N
۳	ابری	"	"	نیست	P
۴	بارانی	معتدل	"	"	P
۵	"	سرد	ترمال	"	P
۶	"	"	"	هست	N
۷	ابری	"	"	"	P
۸	آفتابی	معتدل	بالا	نیست	N
۹	"	سرد	ترمال	"	P
۱۰	بارانی	معتدل	"	"	P
۱۱	آفتابی	"	"	هست	P
۱۲	ابری	"	بالا	"	P
۱۳	"	گرم	ترمال	نیست	P
۱۴	بارانی	معتدل	بالا	هست	N

با استفاده از اطلاعات جدول فوق (مجموعه اجرایی) و بر اساس الگوریتم ID3، درخت تصمیم گیری نمودار ۱ حاصل می شود. که بر اساس درخت حاصل قواعد تصمیم گیری به صورت ذیل خواهد شد.



نمودار ۱ - درخت تصمیم گیری ساده

قواعد تصمیم گیری

۱- اگر هوای صبح روز شنبه آلتاچی و رطوبت بالا باشد آتروز برای فعالیت بخصوص نامساعد است .

۲- اگر هوای صبح روز شنبه آلتاچی و رطوبت نرمال باشد آتروز برای فعالیت بخصوص مساعد است .

۳- اگر هوای صبح روز شنبه ابری باشد آنگاه آتروز برای فعالیت بخصوص مساعد است .

۴- اگر هوای صبح روز شنبه بارانی و طوفانی هم باشد آنگاه آتروز برای فعالیت بخصوص نامساعد است .

۵- اگر هوای صبح روز شنبه بارانی و طوفانی نباشد آنگاه آتروز برای فعالیت بخصوص مساعد است .

در دستورالعملهای ID3 ، بعضی مواقع در مرحله دوم لازم است که خمیمه دارای بیشترین اطلاعات را تعیین کنیم. برای این منظور از تکنیک زیر استفاده می‌کنیم .

در این روش ابتدا اطلاعات مربوط به کل نمونه‌ها را بدست آورده سپس بر اساس آن، میزان موثر بودن تمامی متغیرها را محاسبه می‌کنیم و در میان این مقادیر ، بزرگترین مقدار را مشخص کرده و در نتیجه متغیر مربوطه ، دارای بیشترین اطلاعات می‌باشد .

فرض کنید متغیر کلاس مربوط به مجموعه اجرایی و جدول داده‌ها دارای  $K$  رده باشد. همچنین فرض کنید خمیمه (متغیر)  $A$  دارای مقادیر  $A_1$  و  $A_2$  و ... و  $A_V$  و تعداد نمونه‌های مجموعه (گروه)  $S$  برابر  $n$  باشد. برای محاسبه میزان عدم اطلاع غیر شرطی کلاس (آنتروپی کل کلاس) از فرمول زیر استفاده می‌کنیم .

$$I(C_1, \dots, C_K) = - \sum_{i=1}^K P_{i1} \log P_{i1} \cdot P_i = \frac{\#C_i}{n} = \frac{C_i}{S} \quad \begin{array}{l} \text{تعداد نمونه‌های دارای کلاس } C_i \\ \text{که در آن} \\ \text{تعداد کل نمونه‌ها در } S \end{array}$$

اکنون برای محاسبه میزان موثر بودن متغیر دلخواه  $A$ ، ابتدا متوسط آنتروپی کلاس با فرض متغیر را بدست آورده و از آنتروپی کل متغیر کلاس کم کرده و مقدار حاصل میزان موثر بودن  $A$  برای پیش بینی کلاس خواهد بود. و برای بدست آوردن متوسط آنتروپی کلاس با فرض متغیر  $A$ ، براساس فرمول میانگین موزون زیر عمل می‌کنیم.

$$E(A) = \sum_{i=1}^V \frac{\#A_i}{n} I(C_1, C_2, \dots, C_K / A_i)$$

که در آن  $I(C_1, C_2, \dots, C_K / A_i) = - \sum_{j=1}^K P_{ij} \log P_{ij}$  تعداد نمونه‌هایی که بر اساس افراز  $A_i$  دارای کلاس  $C_j$  هستند

$P_{ij} = \frac{\#(C_j / A_i)}{\#A_i}$  و در آن  $\frac{\#(C_j / A_i)}{\#A_i}$  = تعداد نمونه‌هایی که در افراز  $A_i$  قرار دارند

و در نتیجه :

$$A \text{ میزان موثر بودن} = I(C_1, C_2, \dots, C_K) - E(A)$$

این میزان برای تمامی متغیرها محاسبه و ماکزیمم میزان را مشخص می‌کنیم که سفیر مربوطه، همان متغیر دارای بیشترین اطلاعات خواهد بود. چون در اینجا  $I(C_1, C_2, \dots, C_n)$  برای تمامی متغیرها یکسان است. بنابراین این کافی است. متوسط آنتروپی کلاس با فرض متغیرها را محاسبه و از بین متوسط‌ها، می‌نیم را مشخص کرده و متغیر مربوطه همان متغیر دارای بیشترین اطلاعات خواهد بود. این محاسبات را برای مثال ۱ انجام می‌دهیم. در اینجا تعداد نمونه‌ها  $n = 14$  و تعداد کلاسها  $K = 2$  و بعنوان مثال برای متغیر پیش بینی  $V = 3$  است. ابتدا برای متغیر پیش بینی که دارای سه مقدار آنتابی، ابری و بارانی می‌باشد متوسط آنتروپی را بدست می‌آوریم.

برای مقدار آنتنایی داریم :  $\#N1 = 3$  و  $\#P1 = 2$  و بنابراین

$$I(N, P/A_1 = \text{آنتابی}) = -\frac{9}{14} \log_2 \frac{9}{14} - \frac{5}{14} \log_2 \frac{5}{14} = 1.94 \text{ بیت}$$

وبطور مشابه برای مقادیر ابری و بارانی داریم :

$$I(N, P/A_2 = \text{آبری}) = 0 \quad \#N_2 = 0 \quad \#P_2 = 4$$

$$I(N, P/A_3 = \text{بارانی}) = 1.971 \quad \#P_3 = 3 \quad \#N_3 = 2$$

و در نتیجه :

$$E(\text{پیش بینی}) = \frac{5}{14} I(N, P/A_1 = \text{آنتابی}) + \frac{4}{14} I(N, P/A_2 = \text{آبری}) + \frac{5}{14} I(N, P/A_3 = \text{بارانی}) = 1.694$$

و بنابراین

$$\text{بایت } 0.246 = E(\text{پیش بینی}) - 0.94 = \text{میزان موثر بودن متغیر پیش بینی}$$

وبطور مشابه برای دیگر متغیرها داریم .

$$\text{بایت } 0.29 = \text{موشر بودن دما}$$

$$\text{بایت } 0.151 = \text{موشر بودن رطوبت}$$

$$\text{بایت } 0.48 = \text{موشر بودن طوفانی}$$

و از روی این مقادیر واضح است که متغیر پیش بینی بعنوان متغیر دارای بیشترین اطلاعات تعیین می شود. و نمونه ها بر اساس مقادیر این متغیر تقسیم می شوند و هر افزازی از این متغیر بعنوان یک گروه ( مجموعه S ) در نظر گرفته و با توجه به متغیرهای باقی مانده ، مرحله دوم دستورالعمل های ID3 را روی آن پیاده می کنیم. یکی از حالت های خاص وقتی اتفاق می افتد که یکی از مجموعه های افزاز شده ( گروه ها ) خالی از نمونه باشد. در این حالت ID3 برای آن گروه مقداری ( Null ) را در نظر گرفته و در پایان اجرای الگوریتم ، به گروه های تهی. یکی از کلاس هایی را که بیشترین فراوانی را در مجموعه نمونه ها داشته



باشد، انتخاب می‌دهد .

حالت خاص دیگروقتی اتفاق می‌افتد که یکی از متغیرها دقیقاً "منطبق بر طبقه بندی باشد. فرض کنید متغیر A دارای مقادیر A1 و ۰۰۰ و AV و مجموعه نمونه‌های S را بر اساس مقادیر متغیر A افراز می‌کنیم. چنانچه برای هر یک از مجموعه‌های افراز شده، نسبت نمونه‌هایی را که دارای کلاس C ( $i=1,2,\dots, K$ ) هستند، برابر نسبت نمونه‌ها در کل مجموعه S دارای آن کلاس باشند، آنگاه این متغیر را دقیقاً "منطبق بر طبقه بندی گوئیم. بنابراین وجود این متغیر در درخت تصمیم گیری ضرورتی نخواهد داشت .

الگوریتم ID3 برای حل این مسئله بر اساس میزان موثر بودن متغیرها، از متغیرهای غیر موثر صرف نظر می‌کند. و درخت تصمیم گیری را بر اساس متغیرهای موثر تشکیل می‌دهد. بدینمورت که چنانچه برای هر متغیر، میزان موثر بودن آن از قدر مطلق یا درصدی از یک معیار مشخص کوچکتر باشد، متغیر را غیر موثر شناخته و در محاسبات شرکت نمی‌دهد. در غیر اینمورت موثر شناخته و بر اساس روشی که گفته شد، بدنبال موثرترین متغیر و سپس ادامه الگوریتم می‌باشد .

این الگوریتم برای مقادیر گمشته و اطلاعات دارای نویز ( NOISE ) نیز راه‌های مناسبی ارائه می‌دهد. جهت اطلاع بیشتر در این زمینه همچنین در رابطه با الگوریتم Regressio Tree می‌توانید به پایان نامه تحصیلی اینجانب سید

باقر میراشرقی دانشکده علوم دانشگاه فردوسی مشهد رجوع نمایید .

نرم افزارهایی نیز در این رابطه از سوی شرکت های مختلف تهیه و مورد استفاده قرار می‌گیرند .

دو نمونه از این نرم افزارها تحت عناوین KET , NEWID وده که به همراه راهنمای مربوطه در اختیار قرار دارد. درکنار این دو نرم افزار ، نرم افزار

دیگری نیز اینجانب با همکاری گروه کامپیوتر دانشگاه فردوسی مشهد تهیه  
شده و مورد استفاده قرار دادیم .

#### ۴- نتیجه گیری

همانگونه که در این مقاله آمده است، مابذنبال سیستمی هستیم که با استفاده  
از آن بتوانیم براحتی وبا دقت بیشتر و نیز سرعت زیاد، تجزیه و تحلیلهای  
مختلف آماری را که در زمینه‌های مختلف کاربردی مورد استفاده قرار می‌گیرد،  
به انجام برسانیم .

برای نیل به این هدف از سیستم‌های خبره‌ای که بر اساس روند فکر کردن افراد  
خبره تشکیل شده باشد، استفاده می‌کنیم. و انتظار داریم که در آینده ای نه  
چندان دور برای تمامی تجزیه و تحلیلهای آماری، سیستم‌های خبره‌ای تهیه و مورد  
استفاده علوم و در خدمت بشر قرار گیرد .

در پایان لازم می‌دانم از جناب آقای ازقاسمی استادیار دانشگاه فردوسی مشهد به  
خاطر زحماتی که در رابطه با این مقاله متقبل شدند تشکر و قدردانی نمایم .

مراجع :

۱- پایان نامه کارشناسی ارشد آمار - سیدباقر میراشرقی - دانشکده علوم - دانشگاه فردوسی " مشهد "

2) Breiman, L. et al (1984), Classification and Regression Tree. Wadsworth Belmont.

3) Nakhaeizadeh, G. (1991 a), Inductive Expert Systems and their Application in Statistics. Erscheint in Proceedings of Statsoft 91.

4) Nakhaeizadeh, G. (1991b), Application Of Machine Learning to solving industrial problems. Eyscheint in Proceedings of SOR 91 Trier, Physica Verlag.

5) Quinlan, J.R. (1986), Induction of Decision Trees. Machine Learning. VOL. 1, 81-106.