



## پیش بینی تقاضای دارو به کمک روشی ترکیبی مبتنی بر *SVR* و *KNN* (چالش پیش بینی تقاضای دارو)

سعید جعفری زاده<sup>۱\*</sup>، فاطمه محمد پور<sup>۲</sup> و \*\*احمد رضا نقش نیلچی<sup>۳</sup>

<sup>۱</sup>دانشگاه اصفهان، saeid.jafarizadeh@gmail.com

<sup>۲</sup>دانشگاه پاسارگاد، fatemehmohamadpoor91@gmail.com

<sup>۳</sup>دانشگاه اصفهان، nilchi@eng.ui.ac.ir

\*سرپرست تیم

\*\*استاد راهنمای پروژه

چکیده - این چالش در مورد پیش بینی مقادیر مصرف دارو به صورت روزانه، و به تفکیک هر بیمارستان تا ۱۰ ماه آتی با توجه به داده های سه سال قبل می باشد، روش ارائه شده در این تحقیق ترکیبی از روش *SVR* و *KNN* است. بر اساس نتایج حاصله از پیاده سازی روش های مختلف بر روی داده های ارائه شده در این چالش به این نتیجه رسیدیم که به دلیل کمبود ویژگی های مناسب (تنها یک ویژگی زمان و روز هفته در اختیار می باشد) و همچنین ماهیت مسئله (پیش بینی برای یک دوره زمانی بسیار طولانی) اغلب مدل های پیش بینی از روش های سنتی و ساده گرفته تا روش های عمیق از جمله *LSTM* و *RNN* و *CNN* در اینجا عملکرد مناسبی را از خود نشان نمی دهند. بنابراین انتخاب روشی جهت استخراج ویژگی مناسب (جهت تغییر فضای ویژگی) در این مسئله بسیار حائز اهمیت است لذا نوآوری ارائه شده در این تحقیق تغییر فضای ویژگی و سپس مدل سازی بر اساس ویژگی جدید می باشد. (ذکر این نکته بسیار حائز اهمیت است که در اینجا از داده های خارج از داده های ارائه شده استفاده نشده است)، نتایج نشان دهنده خطای پایین روش ارائه شده نسبت به روش های ذکر شده می باشد.

کلید واژه - *SVR*، *KNN*، *SVR*، *LSTM*، *CNN*، تقاضای دارو، پیش بینی

پیش بینی کند. این اطلاعات می تواند به پزشکان و محققان در تصمیم گیری های بالینی و تحقیقاتی کمک کند [2].

۳. مسئله پیش بینی عوارض جانبی دارو: در این مسئله، هدف پیش بینی عوارض جانبی یک دارو است. با استفاده از داده های بالینی و اطلاعات مربوط به عوارض جانبی، می توان یک مدل پیش بینی برای شناسایی و پیش بینی عوارض جانبی دارو ساخت. این مدل می تواند به پزشکان و سازمان های بهداشتی کمک کند تا در انتخاب و استفاده از داروها، ریسک های مرتبط با عوارض جانبی را در نظر بگیرند و تصمیمات بهتری داشته باشند [3].

۴. مسئله پیش بینی پاسخ به درمان: در این مسئله، هدف پیش بینی پاسخ بیماران به یک درمان خاص است. با استفاده از داده های بالینی، ویژگی های بیماران و اطلاعات مرتبط، می توانیم مدلی بسازیم که بتواند پاسخ بیماران به یک درمان را پیش بینی کند. این اطلاعات می تواند به پزشکان در انتخاب و تنظیم درمان ها و بهبود نتایج بالینی کمک کند [4-5].

در این تحقیق هدف "مسئله پیش بینی تقاضای دارو" در سطح بیمارستان های یک شهر است. در این مسئله تقاضای دارو هر روز ۱۲ بیمارستان مختلف به تفکیک برای یک سال به ما داده شده و از ما پیش بینی مصرف داروی ۱۰ ماه آینده این ۱۲ بیمارستان به تفکیک خواسته شده است. بدیهی است برای حل این مسئله راه حل های مختلفی وجود دارد. که به صورت کلی به دو دسته زیر تقسیم بندی می شوند:

- روش های کلاسیک: از جمله روش های ساده و متداول برای پیش بینی،

### ۱- مقدمه

مصرف داروها در حوزه بهداشت و درمان اهمیت بسیاری دارد و پیش بینی مصرف داروها نیز از اهمیت بالایی برخوردار است. پیش بینی مصرف داروها به ما کمک می کند تا درمان ها را بهتر برنامه ریزی کنیم، عوارض جانبی را کاهش دهیم و هزینه ها را بهینه سازی کنیم.

در مورد پیش بینی مصرف داروها، می توان از رویکردهای متنوعی استفاده کرد. در زیر به برخی از روش ها و مسائل مرتبط با پیش بینی مصرف داروها اشاره می کنم:

۱. مسئله پیش بینی تقاضای دارو: در این مسئله، هدف پیش بینی تقاضای آینده برای یک دارو است. می توان از داده های گذشته مصرف دارو، اطلاعات پزشکی، مشخصات بیماران و سایر متغیرهای مرتبط استفاده کرد تا به یک مدل پیش بینی برای تقاضای آینده دارو برسیم. این مدل می تواند به متخصصان بهداشت و سازمان های تولید و توزیع دارو کمک کند تا توانایی تامین دارو را بهینه سازی کنند [1].

۲. مسئله پیش بینی مصرف دارو در جمعیت خاص: در این مسئله، هدف پیش بینی مصرف یک دارو در یک گروه خاص از جمعیت است. مثلاً می توانیم بررسی کنیم که چگونه مصرف یک دارو در افراد مبتلا به یک بیماری خاص (مانند دیابت) تغییر می کند. با استفاده از داده های بالینی و زیرمجموعه های مرتبط، می توانیم مدلی بسازیم که بتواند مصرف دارو را در این گروه خاص



نرم (Soft Margin Support Vector Machine) محاسبه کرد. این تابع به این معنی است که برخی از نقاط داده می‌توانند از حاشیه خارج شوند و خطایی را ایجاد کنند، اما این خطا با یک پارامتر به نام  $C$  کنترل می‌شود. هر چه مقدار  $C$  بزرگتر باشد، خطای بیشتری مجاز است و هر چه کوچکتر باشد، خطای کمتری مجاز است. این مفهوم به این معنی است که برخی از نقاط داده می‌توانند از حاشیه خارج شوند و خطایی را ایجاد کنند، اما این خطا با یک پارامتر به نام  $C$  کنترل می‌شود. هر چه مقدار  $C$  بزرگتر باشد، خطای بیشتری مجاز است و هر چه کوچکتر باشد، خطای کمتری مجاز است. این روش از یک تابع هسته (Kernel Function) نیز برای ایجاد یک فضای ویژگی جدید استفاده می‌کند که در آن داده‌ها به صورت خطی قابل جداسازی هستند. معمولاً از توابع هسته‌ی گاوسی (Gaussian Kernel)، چندجمله‌ای (Polynomial Kernel)، سیگموئید (Sigmoid Kernel) و خطی (Linear Kernel) استفاده می‌شود در اینجا از تابع خطی به عنوان تابع کرنل استفاده شده است. برای یافتن تابع پیش‌بینی، این روش از یک تابع هزینه (Cost Function) به نام  $\epsilon$ -insensitive loss استفاده می‌کند که در اینجا تابع هزینه. این تابع هزینه به این معنی است که اگر خطای پیش‌بینی کمتر از یک مقدار ثابت به نام  $\epsilon$  باشد، هزینه صفر است و اگر بیشتر از آن باشد، هزینه برابر با اختلاف خطا و  $\epsilon$  است. برای حل این مسئله، از یک روش به نام Lagrange Multipliers استفاده می‌شود که به دنبال بیشینه‌کردن حاشیه بین داده‌ها و تابع پیش‌بینی است. به طور کلی در نهایت، تابع پیش‌بینی به صورت زیر بدست می‌آید:

$$f(x) = \sum_{i=1}^n (\alpha_i - \alpha_i^*) K(x_i, x) + b \quad (1)$$

که در آن  $n$  تعداد داده‌ها،  $\alpha_i$  و  $\alpha_i^*$  ضرایب لاگرانژ،  $K$  تابع هسته (که در اینجا تابع خطی است) و  $b$  یک عدد ثابت است. برای یافتن این ضرایب و عدد ثابت، از یک روش به نام Sequential Minimal Optimization (SMO) استفاده می‌شود که به صورت تکراری، دو ضریب را به روزرسانی می‌کند تا به شرایط KKT برسند. شرایط KKT شرایط لازم و کافی برای بهینه‌سازی یک تابع هزینه با محدودیت‌های خطی یا غیرخطی است.

## ۲-۲- الگوریتم KNN

الگوریتم KNN یا K-Nearest Neighbor یک روش یادگیری ماشین برای دسته‌بندی و پیش‌بینی مقادیر پیوسته است. این روش بر اساس شباهت نقاط داده با یکدیگر کار می‌کند. برای پیش‌بینی برچسب یا مقدار یک نقطه جدید، این روش از  $K$  نزدیک‌ترین همسایه آن در مجموعه داده آموزشی استفاده می‌کند. برای مثال، اگر  $K$  برابر ۳ باشد، برچسب یک نقطه جدید برابر با برچسب بیشترین تکرار در میان سه نزدیک‌ترین همسایه آن خواهد بود. این روش یک روش غیرپارامتریک است که به معنی این است که هیچ فرضی در مورد توزیع داده‌ها نمی‌کند. این روش می‌تواند با داده‌های عددی و رده‌ای کار کند و برای مسائل مختلفی در دسته‌بندی و پیش‌بینی کاربرد دارد. برای فاصله بین نقاط داده، از معیارهای مختلفی مانند فاصله اقلیدسی، فاصله منهتن، فاصله مینکوفسکی و تابع هسته استفاده می‌شود. برای یافتن نزدیک‌ترین همسایه‌ها، از الگوریتم‌های مختلفی مانند الگوریتم حریصانه، الگوریتم KD-Tree و الگوریتم Ball-Tree استفاده می‌شود.

روش‌های خطی می‌باشند. در این روش‌ها، با استفاده از متغیرهای مستقل مانند تاریخ، فصل، مشخصات بیمارستان و سایر متغیرهای مرتبط، می‌توان یک مدل خطی را ساخت و با استفاده از آن، مصرف دارو را برای ۱۰ ماه آینده پیش‌بینی کرد. روش‌هایی مانند رگرسیون خطی، ماشین بردار پشتیبان، درخت تصمیم، مدل‌های فازی و شبکه‌های عصبی ساده برای این منظور مورد استفاده قرار می‌گیرند.

روش‌های مبتنی بر شبکه‌های عمیق: شبکه‌های عصبی عمیق می‌توانند روش‌هایی قدرتمند برای پیش‌بینی باشند. در حوزه یادگیری عمیق، شبکه‌های عصبی بازگشتی (RNNs) مانند LSTM و GRU و شبکه‌های عصبی پیچشی (CNNs) مورد استفاده قرار می‌گیرند. این شبکه‌ها قابلیت استخراج ویژگی‌های پیچیده از داده‌ها را دارند و می‌توانند در پیش‌بینی مصرف داروها و بسیاری از مسائل مرتبط با پیش‌بینی مورد استفاده قرار گیرند. با این حال روش‌های عمیق نیازمند داده‌های بزرگ و سیستم پردازشی بالا می‌باشند.

به طور کلی، برای انتخاب روش مناسب برای پیش‌بینی مصرف داروها، باید به ماهیت داده‌ها، تعداد متغیرها، حجم داده و هدف نهایی توجه کرد. در اینجا تنها داده‌های در دسترس، زمان و مقدار داروی مصرفی در هر روز است لذا حجم داده و تعداد ویژگی‌ها در وهله اول بسیار ناکافی به نظر می‌رسد. لذا در این تحقیق روشی را برای غلبه بر این موضوع پیشنهاد می‌کنیم. در ادامه ابتدا به بیان روش‌های پیش‌زمینه (SVR و KNN) پرداخته و سپس به بیان روش استخراج ویژگی‌های مناسب برای حل چالش مسئله می‌پردازیم، در نهایت الگوریتم کلی ارائه شده برای حل مسئله شرح داده خواهد شد.

## ۲- روش پیشنهادی

به طور خلاصه در روش ارائه شده ابتدا ویژگی تاریخ را به روز در سال، ماه و روز در ماه تقسیم بندی می‌کنیم تا سه ویژگی جدید حاصل شود و سپس به کمک الگوریتمی مشابه الگوریتم KNN، این ویژگی‌ها را به ۱۰ ویژگی جدید از میان داده‌های داده شده نگاشت کرده و فضای ویژگی جدیدی تولید می‌کنیم و به کمک این ویژگی‌ها یک مدل SVR خطی را آموزش می‌دهیم در نهایت و در مرحله تست برای هر داده جدید مشابه همین عمل را انجام داده و پس از محاسبه نگاشت ویژگی داده جدید از مدل آموزش دیده جهت پیش‌بینی مقدار خروجی داده تست مورد نظر استفاده می‌کنیم. برای درک راحت‌تر این مفهوم ابتدا الگوریتم SVR و KNN را شرح داده و سپس نحوه نگاشت ویژگی را به کمک الگوریتم KNN با ذکر جزئیات شرح می‌دهیم.

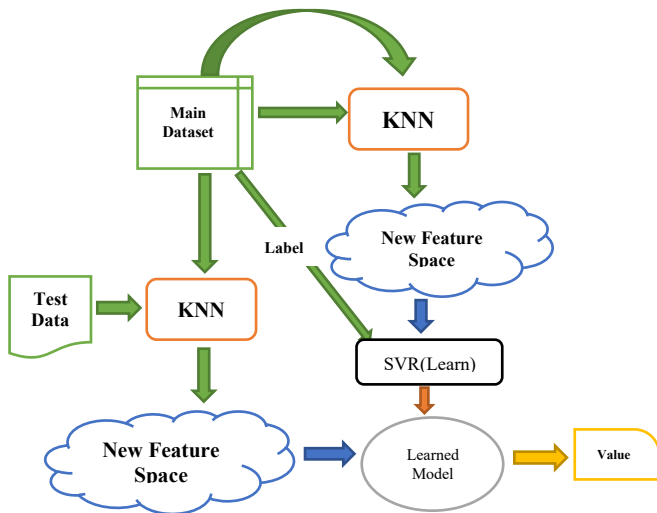
## ۲-۱- الگوریتم SVR

الگوریتم SVR یا Support Vector Regression یک روش یادگیری ماشین برای پیش‌بینی مقادیر پیوسته است. این روش بر اساس ایده‌ی Support Vector Machine (SVM) که یک روش دسته‌بندی است، ساخته شده است. هدف این روش این است که یک تابع خطی یا غیرخطی را بیابد که با حداقل خطای ممکن (MSE)، مقادیر هدف را پیش‌بینی کند. برای این کار، این روش از یک مفهوم به نام ماشین بردار پشتیبان با حاشیه نرم (Soft Margin Support Vector Machine) استفاده می‌کند. در SVR هدف کاهش فاصله بردارهای پشتیبان از خط مارجین می‌باشد. این فاصله را می‌توان با استفاده از یک تابع به نام ماشین بردار پشتیبان با حاشیه

### ۲-۳- نحوه نگاشت ویژگی‌ها به فضای ویژگی جدید به

#### کمک الگوریتم KNN

در اینجا به دلیل ماهیت مسئله که رگرسیون می باشد از روش KNN برای برچسب دهی و دسته بندی استفاده نشده است بلکه تنها برای بدست آوردن داده مشابه به نمونه در حال آموزش یا تست بهره گرفته ایم به این صورت که  $K$  نزدیکترین داده را از بین داده های آموزشی انتخاب می کنیم سپس برچسب این  $K$  داده را به عنوان ویژگی جدید ( و بدون در نظر گرفتن ویژگی های قبلی) برای مرحله آموزش و یا تست مورد استفاده قرار میدهیم، به صورت نشان داده شده در شکل زیر:



شکل ۲: تغییر فضای ویژگی از حالت ویژگی های مبتنی بر زمان به ویژگی های مبتنی بر برچسب خروجی

به طوری کلی در این روش همانطور که در شکل بالا مشاهده می شود از داده های آموزشی هم در مرحله تست و هم در مرحله آموزش استفاده شده است ، در این روش بر خلاف روش های معمول مدلسازی (همانند SVR) پس از آموزش نیز نیازمند داده های دیتاست می باشیم و بر خلاف روش هایی مانند KNN در مرحله آموزش اقدام به ایجاد مدلی خطی از داده و تنظیم پارامتر های مدل مذکور نموده و در مرحله تست نیز از این مدل استفاده مینماییم! در این روش پارامتر  $K$  برابر با ۱۰ در نظر گرفته شده و کرنل SVR نیز همانطور که گفته شد به صورت خطی در نظر گرفته شده است (نتایج نشان دهنده عملکرد بهتر کرنل خطی نسبت به کرنل های غیر خطی در این مسئله است.) همچنین مدل SVR مورد نظر در زمان آموزش به کمک روش اعتبار سنجی 10-fold بهینه می شود .

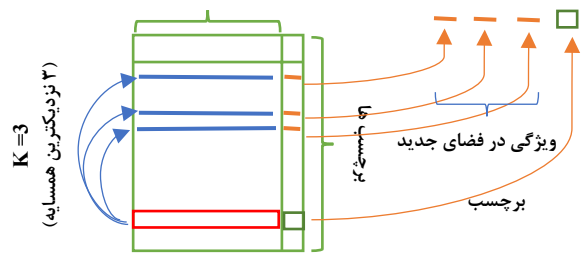
### ۳- نتیجه گیری

در این تحقیق علاوه بر روش پیشنهادی از روش هایی همانند LSTM دو گانه و یکتا به صورت ارائه شده در مرجع [6] و CNN (به صورت یک بعدی بر روی پنجره ای از داده ها به عنوان سری زمانی) نیز جهت حل مسئله استفاده شده است. اما نتایج نشان دهنده عملکرد ضعیف تر این الگوریتم ها در قبال روش ترکیبی ارائه شده می باشد . همچنین با توجه به تغییر فضای ویژگی ، ویژگی های تغییر یافته را نیز برای هر دو روش LSTM و CNN گفته شده بکار گرفته ایم اما کماکان نتایج حاصله نشان دهنده دقت بالاتر روش پیشنهادی می باشد بطوری که اختلاف خطای ۵ تا ۱۰ برابری نسبت به روش پیشنهادی در زمان اعتبار سنجی را شاهد بودیم. برای مثال جدول ۱ نشان دهنده بهترین نتایج حاصله برای پیش بینی داده های بیمارستان ۱ می باشد

جدول ۱ - مقایسه روش ارائه شده با مدل های مشابه برای بیمارستان ۱

روش	CNN	LSTM	SVR	Our Method
معیار				
MAE	17.60	7.05	17.61	0.7601

#### ویژگی ها



شکل ۱: نحوه تغییر فضای ویژگی با توجه به داده های موجود در مسئله پیش بینی تقاضای دارو

به عبارت دیگر در اینجا فضای ویژگی را از فضای زمانی به فضای برچسب ها نگاشت کرده ایم ، در ادامه می بینیم با توجه به نتایج حاصله در این فضا مدلی هایی مانند بردار پشتیبان ماشین عملکرد بهتری را نسبت به روش های دیگر از خود نشان می دهند.

### ۲-۴- مدلسازی و آموزش

حال که روش تغییر فضای ویژگی ها شرح داده شد، در ادامه از روش SVR جهت مدلسازی و پیش بینی داده ها استفاده می کنیم، ابتدا هر یک از داده های دیتاست را به فضای جدید برده و سپس داده ها را بدون اعمال نرمال سازی به مدل SVR خطی می دهیم و مدل را با آن آموزش می دهیم ، در این مرحله برای هر یک از داده ها مقدار پیش بینی را از دیتاست انتخاب می کنیم. به عبارت دیگر برچسب مربوطه را به ازای هر یک از داده ها از دیتاست مورد نظر بیرون کشیده و مدل مورد نظر را آموزش می دهیم در مرحله تست نیز ابتدا ویژگی های مورد نظر را در فضای جدید از پایگاه دانش به کمک الگوریتم KNN داده شده بدست می آوریم و در نهایت مقدار پیش بینی نهایی را به کمک مدل آموزش داده شده در مرحله قبل محاسبه می کنیم ، شکل ۲ نشان دهنده نحوه عملکرد روش پیشنهادی مورد نظر و نحوه پیش بینی داده های جدید می باشد.



راستای انجام این پروژه یاری رسان بودند.

### منابع

- [۱] بیگلرخانی امین، عباسی رضوان، ثنایی محمدرضا. "پیش‌بینی میزان مصرف دارو در بیمارستان‌ها با استفاده از مدل شبکه حافظه طولانی کوتاه‌مدت. بیمارستان". ۱۴۰۱؛ ۲۱ (۴): ۲۲-۳۵.
- [۲] مهاجر تبریزی، خوجه، درویش محمدی. "طراحی شبکه توزیع دارو با استفاده از الگوریتم ژنتیک دوسطحی و کدینگ مبتنی بر اولویت (مطالعه موردی: شرکت پخش سراسری آدوراطب)". مدیریت تولید و عملیات، ۲۰۲۲، ۱۳(۳)، ۴۷-۷۵.
- [۳] سندسی، گلشن، صائی، هاشمی گلپایگانی. "یک روش پیش‌بینی پیوند مبتنی بر همسایه برای شبکه دویخشی" فصلنامه فناوری اطلاعات و ارتباطات ایران، ۴۷(۴۷)، ۱۷۵.
- [۴] رستمی. پیش‌بینی کننده‌های بالینی پاسخ به درمان تحریک مکرر مغناطیسی فراجمجمه ای (rTMS) در بیماران مبتلا به اختلال افسردگی. فصل‌نامه پژوهش‌های کاربردی روانشناختی، ۲۰۱۷، ۲: ۸۰-۱۸۱.
- [۵] گلبائنی سروش، برهانی خاطره برهانی حامد. پیش‌بینی تصمیم‌گیری اخلاقی کادر درمان در شرایط غیرقطعی بر اساس مواجهه با مرگ، احتمال ابتلا و رضایت شغلی: نقش میانجی‌گری اضطراب. پژوهش‌های روانشناسی اجتماعی ۲۰۲۳، ۱۳(۴۹)، ۱-۱۴.
- [6] GAJAMANNAGE, Kelum; PARK, Yonggi; JAYATHILAKE, Dilhani I. "Real-time forecasting of time series in financial markets using sequentially trained dual-LSTMs". Expert Systems with Applications, 2023, 223: 119879.

در اینجا معیار ارزیابی میانگین مطلق خطا (MAE) در نظر گرفته شد است که به صورت فرمول ۳ محاسبه می‌شود.

$$MAE = \frac{1}{n} \sum_{i=1}^n |\hat{y}_i - y_i| \quad (3)$$

که در آن  $\hat{y}$  مقدار پیش‌بینی شده و  $y$  مقدار واقعی میزان مصرف دارو در روز می‌باشد. به عنوان تحلیل نتایج حاصله می‌توان دلیل عدم موفقیت مدل‌های عمیق را در این مسئله به دلیل نداشتن داده کافی، عدم غنای کافی فضای ویژگی‌ها، انتظار پیش‌بینی طولانی مدت (۱۰ ماه آینده) و مهمتر از همه تغییر شدید در داده‌های آموزشی بیان کرد به طوری که ممکن است در یک بیمارستان میزان مصرف دارو در یک روز رقم بالا (مثلاً ۶۰ قلم دارو) و فردای آن روز تعداد بسیار پایین (مثلاً ۱ یا دو قلم دارو) باشد (این امر هیچ وابستگی به ویژگی‌های ارائه شده مانند تعطیلی روز قبل یا بعد نیز ندارد و به کرات در داده‌دیتاست مشاهده می‌شود و از آنجایی که این امر به تعداد بسیار زیاد در دیتاست مشاهده شده نمیتوان این داده‌ها را به عنوان داده پرت در نظر گرفت!). طبیعتاً با این اوصاف، مدل‌هایی همانند LSTM توانایی پیش‌بینی تغییرات ناگهانی و پیش‌بینی دراز مدت را نداشته و لذا در اینجا بدیهی است که پیش‌بینی درستی را برای این تغییرات ارائه ندهند همچنین مدل‌هایی مانند CNN نیز بیش از حد به رفتار محلی سری زمانی وابستگی داشته و به تاریخ‌های مشابه همانند تاریخ‌های تکراری در سال یا ماه توجهی ندارد لذا روش پیشنهادی به نظر از لحاظ منطقی نیز بایستی به دلیل عدم وابستگی زیاد به حافظه کوتاه یا بلند مدت و همچنین عدم توجه زیاد به داده‌های محلی در آموزش نتیجه بهتری نسبت به مدل‌های مشابه داشته است.

### سپاسگزاری

با تشکر از اساتید گرانقدر گروه هوش مصنوعی دانشگاه اصفهان که ما را در