

# HERFEH+

حرفه پلاس، ورود به دنیای حرفه‌ای‌ها

## داده کاوی

## تعريف داده کاوی:

زمانی که از بحث داده کاوی صحبت می‌شود دو مبحث را باید در نظر گرفت یکی بحث داده یا Data و دیگری بحث اطلاعات یا Information گفته می‌شود که اطلاعات قدرت می‌آورد چون می‌تواند تصمیمات به هنگام انجام دهد ولی موضوع این است هر چیزی که در گام اول دریافت می‌شود قابل تصمیم‌گیری نیست ما ممکن است داده‌های زیادی دریافت کنیم اما این داده‌ها ممکن است داده‌های درستی نباشند یا حجم این داده‌ها زیاد باشد یا تنوع این داده‌ها زیاد باشد یا مسائل دیگری وجود داشته باشد. وقتی ما می‌توانیم تصمیم‌گیری درستی برای کسب‌وکار خود بگیریم که این داده‌ها را به نوعی کاویده باشیم و آن چیزی که برای کسب‌وکار ما ارزشمند است برای تصمیمات استفاده کنیم به عبارتی ما باید با فیلتر کردن داده‌هایی که به دست ما می‌رسد بخشی که خیلی ارزشمند در جهت منافع سازمانی ما هست را پیدا کنیم و آن‌ها را به اطلاعات تبدیل کنیم و از آن‌ها استفاده کنیم به عبارت دیگر داده‌ها یک مبحث خام است و اطلاعات یک مبحث پخته و کامل است که عملاً از دیتای اولیه به دست آمده است.

داده کاوی در حقیقت استخراج اطلاعات از داده‌های اولیه است بنابراین هوشمندی بسیاری می‌خواهد تا بتوانیم از داده‌های اولیه اطلاعات موردنیاز خود را به دست آوریم.

وظایف داده کاوی یا data mining چیست؟

کاویدن و استخراج دانش از منابع عظیم داده برای آشکارسازی اطلاعات گران‌بهای پنهان شده در حجم انبوه اطلاعات سطحی. مثلاً کجا سرمایه‌گذاری کنیم، یا کجا محصولمان را بفروشیم یا کجا خاط تولیدمان را بیشتر کنیم یا کجا وارد بورس شویم یا نشویم. و چشم بسته در زمینه حرفه خود تصمیمات مدیریتی نگیریم.

## پایه‌های داده کاوی شامل:

- آمار (مطالعه عددی روابط داده‌ها)
- هوش مصنوعی (هوش انسان مانند نرم‌افزار و یا ماشین)
- یادگیری ماشین (الگوریتم‌هایی برای پیش‌بینی داده‌ها)

## تفاوت داده کاوی و تحلیل‌های علم آمار

داده کاوی در حقیقت زیرمجموعه علم آمار است اما تفاوت‌هایی نیز بین آن‌ها وجود دارد که شامل:

- حجم داده‌های مورد تحلیل
- روش مدل‌سازی داده‌ها
- استفاده از هوش مصنوعی
- فراتر رفتن اهمیت داده کاوی از علم آمار
- در اختیار قراردادن امکانات
- امکان استفاده از مدل‌های مدل‌سازی داده‌های ورودی
- دستیابی به اطلاعات پنهان موجود در حجم عظیم داده
- استفاده از شیوه‌های خودکار یادگیری

## ویژگی‌ها و فواید داده کاوی یا DM

- ویژگی‌های اصلی داده کاوی عبارت‌اند از:
- کشف اتوماتیک الگوهای در انجام پروژه‌ها
- پیش‌بینی احتمالی نتایج و خروجی‌ها
- تولید اطلاعات اجرایی و مفید
- مرکز بر داده‌های بزرگ
- ایجاد روابط به صورت اتوماتیک
- استفاده از داده‌های متنوع
- داینامیک بودن
- عدم نیاز به داده‌های صحیح
- ساخت مدل‌های واقعی
- آنالیز کردن داده‌های واقعی
- دوری از اشکالات احتمالی نمونه‌گیری

## برخی از فایده‌های داده کاوی:

- تصمیم‌گیری‌های واقع‌بینانه
- کاهش هزینه‌ها
- تکرار تصمیمات سودآور
- بالابردن میزان فروش
- ارتقای کیفی پژوهش‌ها
- ارزیابی میزان ریسک
- حذف ضررهای ناشی از ناآگاهی مدیران
- شفاف کردن فضای حاکم برکسب وکار
- تشخیص جرائم و شناسایی سوابق افراد

## برخی از محدودیت‌های داده کاوی:

- مشخص نبودن ارزش یا میزان اهمیت الگوها و روابط
- نیازمند بودن به کارشناسان فنی و تحلیل گران کارآزموده
- نیاز به صرف هزینه‌های بالا در برخی موارد

## نمونه کاربردهای داده کاوی

- استفاده تجاری: مانند الگوی خرید افراد از فروشگاه‌های زنجیره‌ای
- استفاده پزشکی: مانند کشف الگوهای ناشناخته تاثیر داروها بر بیمارهای مختلف و گروههای سنی مختلف
- استفاده در بانک داری: سرمایه‌گذاری و بورس مانند شناسایی مشتریان پرخطر یا سودجو

## چگونگی عملکرد سیستم و فرآیندهای داده کاوی

گام اول: درک کسب‌وکار

گام دوم: بررسی و درک داده‌ها

گام سوم: آماده‌سازی داده‌ها

گام چهارم: مدل‌سازی

گام پنجم: تست و ارزیابی مدل

گام ششم: توسعه مدل نهایی

سه گام اول در حقیقت گام‌های مقدماتی و سه گام دوم گام‌های اجرایی است.

**۱-درک کسب‌وکار:** در میان گذاشتن مسئله کسب‌وکار با متخصص داده کاوی و واگذاری تشخیص مسئله کسب‌وکار به متخصص داده کاوی.

**۲-بررسی و درک داده‌ها:** بررسی داده‌ها با توجه به حجم و کیفیت داده‌ها و تعديل مسئله برای واقع‌بینانه شدن پروسه داده کاوی.

شدن پروسه داده کاوی.

۳- آماده سازی داده ها: یکی کردن انبارهای داده در کسب و کار و شناسایی و حذف داده های اشتباه همچنین تغییر فرمت داده ها با توجه به مسئله تعديل شده در گام دوم فرآیند

۴- مدل سازی: انتخاب مدل با توجه به متدهای مختلف شامل:

- مدل سازی توصیفی

- مدل سازی پیش بینانه

- مدل سازی تجویزی

- مدل سازی وابستگی

است که باید بر اساس شرایط کسب و کار تعریف شده باشد.

۵- تست و ارزیابی مدل: تست و ارزیابی مدل انتخاب شده و بررسی موثر بودن و میزان تاثیر و در صورت لزوم تکرار این تست و ارزیابی

۶- توسعه مدل نهایی: ارائه راه حل ها با توجه به مسئله و ارائه نرم افزار شبیه سازی کسب و کار

## مدل استاندارد برای داده کاوی کسب و کار

در بحث پروسس مدل دیگری را در موردش گفتگو می کنیم به نام Cross-Industry Standard Process CRISP-DM که مخفف جمله for Data Mining است که با دو تفکر به موضوع نگاه می کند که یکی در مورد سگمنتیشن است که می گوید من داده های اولیه ای که دارم پس از اینکه گام ها و فرآیندهای داده کاوی که در بالا ذکر شد را انجام داد همه این مراحل می تواند بر اساس دسته بندی یا Classification می تواند انجام شود و می توانیم آنها را به صورت خوش بندی یا Clustering کنیم و بعد به سراغ تحلیل قواعد انجمنی یا Association Rules برویم و سیستم های مصور سازی یا Visualization را اجرا کنیم البته اینها به صورت AND , OR است یعنی همه را می توانیم داشته باشیم و یا بعضی از مکانیسم ها را.

در کنار این مسئله می‌تواند داده کاوی با نظارت ماکه به آن یادگیری با ناظر یا Supervised Learning یا یادگیری بدون ناظر یا Unsupervised Learning یا به صورت سیستم‌های میانگین یا یادگیری نیمه ناظر Semi-Supervised Learning انجام شود.

## معرفی ابزارهای داده کاوی:

ابزارهای داده کاوی خیلی زیاد هستند برخی از نرم‌افزارها بسیار سنگین هستند برخی سبک و اتفاقاً کارا هستند ما باید از کدامیک از ابزارها استفاده کنیم، آیا می‌توانیم از چند ابزار در کنار هم استفاده کنیم. جواب بله است.

بهتر است که از یک مدل استفاده کنیم اما طبیعتاً می‌توانیم از چند مکانیسم استفاده کنیم.

مثلاً برخی از داده‌ها که نیاز به خوشبندی است می‌توانیم سراغ نرم‌افزارهایی برویم که کارشان خوشبندی است، برخی با آمار و ارقام کار می‌کنند برخی با خصلتها کار می‌کنند برخی از زبانهای برنامه‌نویسی با برنامه‌نویسی دیگری نوشته شده‌اند مثلاً پایه‌شان پایتون یا C++ باشد

بنابراین بسته به تفکری که پشت این داده کاوی است در کوتاه مدت می‌تواند تاثیرات شگرفی را ایجاد کند.

در اینجا برخی از نرم افزارهایی که بسیار کاربری هستند و قابل دسترسی و به راحتی قابل بارگذاری هستند به شما معرفی می‌کنیم:

## Rapid Miner

این نرم افزار دارای طیف گسترده‌ای از برنامه‌های کاربری است، می‌تواند برای پروژه‌های تجاری، علم و تحقیق دانشگاه، آموزش و غیره مورد استفاده قرار گیرد. کاربران می‌توانند با استفاده از این ابزار از تکنیک‌های مختلف داده کاوی استفاده کنند.

اگرچه این یک نرم افزار جدید نیست، هنوز هم به عنوان یکی از بهترین‌های شناخته شده است و در مورد قیمت‌گذاری خیلی انعطاف‌پذیر است.

## Orange

این ابزار داده کاوی از بهترین نرم افزارهای متن بازی است که برای مصورسازی داده روی وب موجود است. این نرم افزار به پردازش متن و یادگیری ماشین اختصاص دارد. همچنین دارای برخی از ویژگی‌های پیشرفته برای تجزیه و تحلیل است. با اینکه نرم افزار ساده‌ای است اما می‌تواند کارهای زیادی برای شرکت شما انجام دهد.

## GraphLab Create

این نرم افزار، بهینه‌سازی مدل، رگرسیون، مصورسازی یادگیری ماشین، تشخیص آنومالی، ساختار داده‌های مقیاس‌پذیر، تجزیه و تحلیل تصویر، یادگیری داده‌ها را ارائه می‌دهد. این ابزار اجازه می‌دهد تا شما در مورد مجموعه داده‌ها بیشتر بدانید و الگوهای مختلف را با استفاده از تکنیک‌های مختلف پیدا کنید.

با این کار، می‌توانید برنامه‌های پیش‌بینی کننده بسازید. چنین ابزاری برای شرکت‌ها شگفت‌انگیز است زیرا دوره‌های برنامه‌ریزی آن‌ها را کاهش می‌دهد و توسعه پروژه‌ها را تسريع می‌بخشد.

## R Studio

آر یکی از ابزارهای داده کاوی رایگان در این لیست است. بسیاری از شرکت‌ها با محاسبات آماری و گرافیک در حال مبارزه هستند. این ابزار به شما این امکان را می‌دهد که یاد بگیرید که این ابزار همه‌چیزهایی است که شما نیاز دارید تا دانش خود را در مورد موضوع افزایش دهید.

محیط نرم‌افزار آر اجازه می‌دهد تا شما محاسبات گرافیکی مختلف انجام دهید. این ابزار برای مدل‌های خطی و غیرخطی، طبقه‌بندی، خوش‌بندی و غیره مناسب است. شما می‌توانید با زبان اسکریپتی مختلف به آن دسترسی پیدا کنید.

## Weka

هنگامی‌که از طریق منوی ویکا جستجو می‌کنید، به سرعت متوجه می‌شوید که این نرم‌افزار بر اساس چندین الگوریتم یادگیری ماشین استوار است. برخی از ویژگی‌های ویکا شامل پردازش، طبقه‌بندی، رگرسیون، ارتباط، انتخاب ویژگی، آزمایش‌های مختلف و غیره است. ویکا از اکسپلور استفاده می‌کند.

## KNIME

این نرم‌افزار پلتفرمی است که بر سه چیز تمرکز دارد: ادغام، تجزیه و تحلیل و گزارش دهی. با استفاده از این ابزار داده کاوی، کارهای زیادی می‌توانید انجام دهید. بیش از ۱۰۰۰ تجزیه و تحلیل مختلف وجود دارد که می‌توانید با آن انجام دهید. این ابزار برای تحلیل خود از منابع مختلف استفاده می‌کند. اگر تصمیم به اضافه کردن آن به لیست نرم‌افزاری خود داشته باشید قطعاً ارزش آن را دارد.

این نرم افزار عمدتاً بر روی پیاده سازی شبکه های عصبی تمرکز دارد. این پلت فرم در سی پلاس پلاس ساخته شده و مزیت اصلی آن کارایی بالا است. با توجه به مدیریت حافظه و سرعت پردازش بالا، به راحتی می تواند نرم افزارهای رقیب را مغلوب کند.

## A p a c h e

نام این نرم افزار به معنای برنامه مدیریت اطلاعات بدون ساختار است. به عبارت دیگر، مانند بسیاری دیگر از ابزارهای داده کاوی، این ابزار یک الگو را در مجموعه داده های بزرگ دنبال می کند. با آپاچی می توانید برنامه ها را به اجزای اصلی تقسیم کنید. همچنین می تواند این اجزا را به خدمات شبکه متصل کند. اگرچه این ابزار نسبتاً قدیمی است اما به طور مرتب بروز می شود که بدین معنی است که شما بهترین محصول را خواهید گرفت

## CLUTO

همان طور که از نامش برمی آید، این نرم افزار بر خوش بندی به عنوان تکنیک داده کاوی تمرکز دارد. از مجموعه داده های مختلف صرف نظر از اندازه آن ها برای تجزیه و تحلیل خوش بندی استفاده می کند. این نرم افزار شامل هر دو برنامه مستقل و همچنین کتابخانه است که مورد استفاده قرار می گیرد که بر مبنای آن برنامه های کاربردی می توانند به الگوریتم های خوش بندی دسترسی داشته باشند. این ابزار از چنین روش برای خوش بندی به طور صحیح و کارآمد خلاصه می کند.

## Anaconda

یک پلتفرم کامل داده کاوی است که با پایتون ساخته شده است. مانند بسیاری از این برنامه‌ها، مواردی از قبیل ساده‌سازی جریان‌های کاری از نقطه اولیه تا استقرار وجود دارد. همچنین به شما اجازه می‌دهد که منابع خود را بهمنظور به دست آوردن بهترین نتایج و الگوهای ادغام کنید. این نرمافزار برای تیمهای بزرگ ساخته شده است. شما می‌توانید نتایج و پروژه‌های خود را با کل شرکت به اشتراک بگذارید.

## Shogun

این ابزار می‌تواند توسط هرکسی استفاده شود. این نرمافزار بر اساس الگوریتم‌های مختلفی است که ما را قادر به یافتن الگوهای مختلف می‌کند. این نرمافزار در سی‌پلاس پلاس ساخته شده است و از سایر زبان‌های برنامه‌نویسی نیز پشتیبانی می‌کند. شاید یکی از بهترین موارد در مورد این ابزار رایگان بودن آن است.

## TraMineR

یک نرمافزار داده کاوی است که برای موارد مختلفی استفاده می‌شود، با این ابزار شما می‌توانید داده‌ها را تحلیل و مستندسازی کنید.

## ROSETTA

شما می‌توانید از آن برای مرور و پیش‌پردازش اطلاعات در فازهای اولیه داده کاوی، برای اعتبارسنجی و تجزیه و تحلیل قواعد و الگوهای القا شده استفاده کنید. این نرمافزار برای یک برنامه خاص ساخته نشده بلکه با هدف ابزار عام و چندمنظوره ساخته شده است.