

گروه کامپیوتر – دانشکده مهندسی – دانشگاه فردوسی مشهد

## مستندات گزارش بلند

برای بخشی از فعالیت های کلاسی

درس روش پژوهش و شیوه ارائه مطالب علمی و فنی

موضوع:

علوم داده در کسب و کار

Data science in Business

ارائه دهنده:

حمید رضا زهتاب

استاد درس:

دکتر رضا منصفی

زمستان 1400



## چکیده

امروزه بیش از هر زمان دیگری داده داریم و با توجه به سریع تر شدن و قدرتمند شدن کامپیوتر ها و همچنین به لطف تکنولوژی ذخیره سازی اطلاعات می توانیم هزاران هزار ترابایت داده را مهار کنیم . به زبان ساده علم داده به ما کمک می کند تا این حجم گسترده از اطلاعات را برای خود قابل فهم کنیم . با کمک علم داده ، اطلاعات را طبقه بندی و دسته بندی میکنیم سپس آنها را مصور سازی می کنیم و سپس از داده ها نتیجه گیری می کنیم . نتیجه گیری نهایی می تواند در گستره بزرگی از شاخه ها من جمله کسب و کار به ما کمک کند . در این مقاله تلاش می کنیم تا ابتدا درک درستی از آنچه در علوم داده اتفاق می افتد ارایه کنیم و سپس به معرفی کاربرد های علوم داده در کسب و کار و نحوه به کار گیری آنها می پردازیم .

## فهرست مطالب

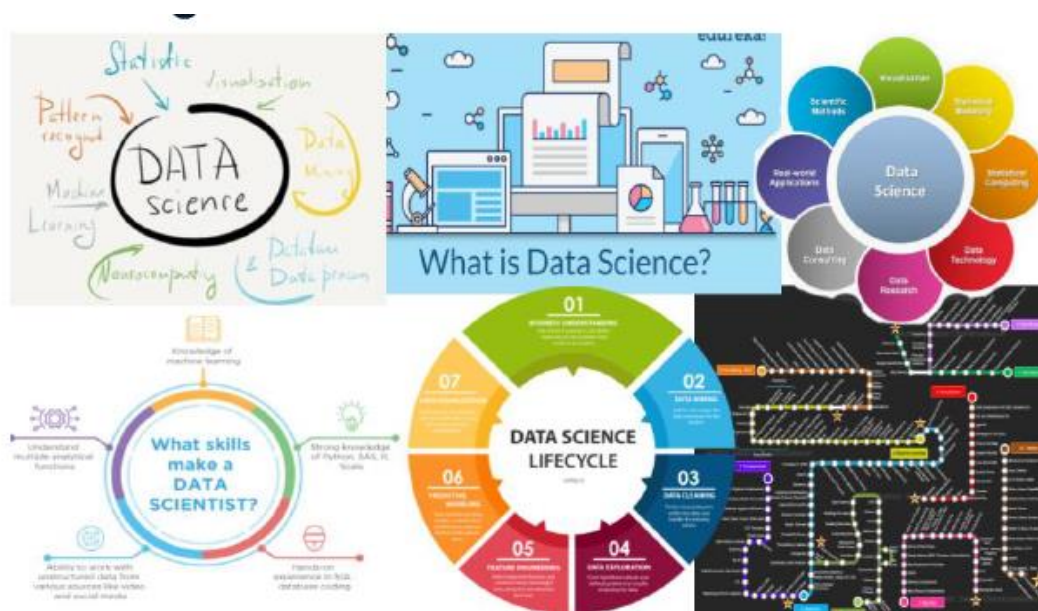
۴	۱-مقدمه .....
۷	۲- کاربرد های علوم داده.....
۷	۲-۱- یادگیری ماشین ( machine learning).....
۸	۲-۲- اینترنت اشياء ( IOT ) .....
۹	۲-۳- یادگیری عمیق.....
۱۱	۳- ساخت تیم علوم داده.....
۱۱	۳-۱- اعضاء تیم.....
۱۱	۳-۱-۱- مهندس داده.....
۱۲	۳-۱-۲- تحلیلگر داده.....
۱۳	۳-۱-۳- دانشمند یادگیری ماشین.....
۱۴	۳-۱-۴- مرور.....
۱۴	۳-۲- ساختار تیم علم داده.....
۱۴	۳-۲-۱- ایزوله شده ( isolated ) .....
۱۵	۳-۲-۲- تعبیه شده ( embedded ) .....
۱۶	۳-۲-۳- ترکیبی ( hybrid ) .....
۱۷	۴- منابع داده و خطرات .....
۱۷	۴-۱- منابع مشترک داده ها.....
۱۷	۴-۲- داده های وب.....
۱۸	۴-۳- اطلاعات قابل شناسایی شخصی ( PII ) .....
۱۸	۴-۴- نام مستعار داده ها.....
۱۹	۴-۵- ناشناس سازی داده ها.....

۱۹.....	۴-۶- مقررات عمومی حفاظت از داده ها.....
۲۰.....	۵- داده های درخواستی.....
۲۰.....	۵-۱- چرا داده ها را درخواست می کنیم؟.....
۲۰.....	۵-۲- انواع داده ای درخواستی؟.....
۲۱.....	۵-۳- ترجیحات آشکار و بیان شده.....
۲۲.....	۵-۴- بهترین شیوه ها.....
۲۳.....	۶- جمع آوری داده های اضافی.....
۲۳.....	۶-۱- حتی داده های بیشتر.....
۲۳.....	۶-۱-۱- API های داده.....
۲۴.....	6-1-2- سوابق عمومی.....
۲۵.....	6-1-3- ترک مکانیکی.....
۲۷.....	7- ذخیره و بازیابی داده ها.....
۲۷.....	7-1- راه حل های ذخیره سازی موازی.....
۲۷.....	7-2- ابر.....
۲۸.....	7-3- انواع ذخیره سازی داده ها.....
۲۹.....	7-4- پرس و جو داده ها.....
۳۰.....	7-5- کنار هم قرار دادن همه : مکان.....
۳۰.....	7-6- کنار هم قرار دادن همه : نوع داده ها.....
۳۱.....	7-7- کنار هم قرار دادن همه : پرس و جو.....
۳۳.....	منابع.....

## ۱- مقدمه

علم داده یک رشته پویا است. ابزار هایی که ما استفاده می کنیم و توانایی های تیم هایمان هر روز در حال تغییر است. در ادامه خواهید دانست که علم داده چیست و چگونه می توانید از آن برای تقویت سازمان خود استفاده کنید.

اگر عبارت "علوم داده چیست؟" را گوگل کنیم، با حجم انبوهی از اطلاعات سردرگم کننده روبرو می شویم، اما علم داده در حقیقت بسیار ساده است. علم داده مجموعه ای از روش ها برای گرفتن هزاران شکل از داده هایی است که امروزه در دسترس ما هستند و از آنها برای نتیجه گیری معنادار استفاده می کنیم. داده ها در همه جا هستند. هر لایک، کلیک، ایمیل، تراکنش بانکی یا توییت یک داده جدید است که می تواند برای توصیف بهتر زمان حال یا پیش بینی بهتر آینده استفاده شود.

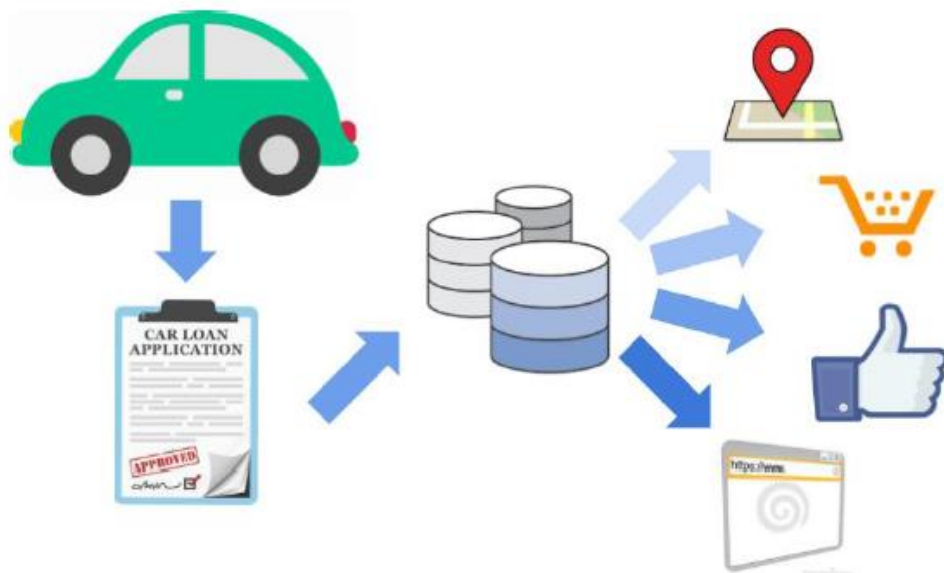


سوال اینجاست که داده ها چه کاری می توانند برای ما انجام دهند؟ داده ها می توانند وضعیت فعلی ما را توصیف کنند. این را می توان با داشبوردها یا هشدارها انجام داد و فرآیندهای زمانبر گزارشدهی را با فناوری های جدید ساده کرد. علم داده می تواند به تشخیص رویدادهای غیر عادی کمک کند. اگر داده هایی در مورد آنچه قبلاً اتفاق افتاده است داشته باشیم، می توانیم با شناسایی خودکار یک رویداد جدید که غیرمنتظره است، کارایی را افزایش دهیم. داده ها همچنین می توانند علل رویدادها و رفتارهای مشاهده شده را تشخیص دهند. به جای تعیین همبستگی بین تعداد کمی از رویدادها، تکنیک های علم داده به ما کمک می کنند تا سیستم های پیچیده را با دلایل احتمالی متعدد درک کنیم. در نهایت، داده ها می توانند رویدادهای آینده را پیش بینی

کنند. ما می توانیم از تکنیک های جدید برای در نظر گرفتن علل مختلف و پیش بینی نتایج بالقوه استفاده کنیم. علاوه بر این، ما می توانیم احتمال پیش بینی خود را به صورت ریاضی ارزیابی کنیم تا سطح عدم قطعیت خود را روشن کنیم.



بنابراین اکنون می دانیم علم داده چیست. سوال بعدی این است که چرا اینقدر محبوب است؟ پاسخ بسیار ساده است: ما در حال جمع آوری داده های بیشتری نسبت به قبل هستیم. فرض کنید به یک نمایندگی خودرو مراجعه می کنید و اطلاعاتی را تکمیل می کنید. همه این داده ها به طور خودکار در رایانه وارد می شوند و با داده های صدها نمایندگی در یک پایگاه داده بزرگ ترکیب می شوند. هنگامی که این داده ها را در اختیار داریم، استفاده از آدرس ایمیلی که هنگام خرید آن خودرو ارائه کرده اید برای پیوند دادن داده های خرید خودرو با داده های شما از رسانه های اجتماعی یا مرور وب آسان است. ناگهان، ما یک تصویر بسیار کامل در مورد همه کسانی که در سال گذشته ماشین خریداری کرده اند داریم: سن آنها، علایق آنها، دوستان و خانواده آنها. این داده های اضافی را می توان برای پیش بینی اینکه چه قیمتی می توانید برای ماشین خود بپردازید، چه خریدهای دیگری انجام دهید، یا بهترین نحوه فروش بیمه آن ماشین جدید به شما، استفاده می شود. داده ها در همه جا وجود دارد و برای کسب و کارها بسیار ارزشمند است.

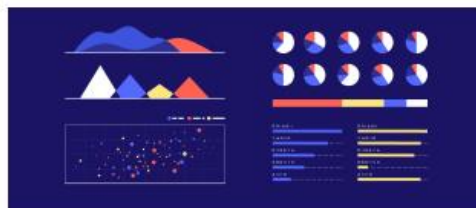


بنابراین، چگونه شروع به استفاده از داده های خود کنیم؟ در علم داده، ما به طور کلی سه مرحله برای هر پروژه داریم. اول، ما داده ها را از منابع بسیاری، مانند نظرسنجی های مشتریان، نتایج ترافیک وب، ایمیل های بین نمایندگان فروش و مشتریان بالقوه، و تراکنش های مالی جمع آوری می کنیم. در مرحله بعد، داده های خود را بررسی و مصور (visualize) می کنیم. این می تواند شامل ساخت داشبوردهایی برای ردیابی اینکه چگونه داده های ما در طول زمان تغییر می کنند یا انجام مقایسه بین دو مجموعه داده باشد. در نهایت، ما با داده های خود پیش بینی می کنیم. به عنوان مثال، این می تواند شامل ساختن سیستمی باشد که مشتریان را تقسیم بندی می کند یا تصاویر انواع خودروها را طبقه بندی می کند.

#### Data collection



#### Exploration and visualization



#### Experimentation and prediction





## 2- کاربرد های علوم داده

قبلاً با تعریف علم داده و مراحل یک گردش کار علم داده آشنا شدید. در ادامه ، شما یاد خواهید گرفت که چگونه علم داده را برای مشکلات واقعی کسب و کار به کار ببرید.

بیایید به سه حوزه مهیج علم داده نگاهی بیندازیم: یادگیری ماشین ، اینترنت اشیا و یادگیری عمیق.

### 2-1- یادگیری ماشین

فرض کنید در یک بانک بزرگ در اداره کشف کلاهبرداری و تقلب ( fraud detection ) کار می کنید. می خواهید از داده ها برای تعیین احتمال جعلی بودن تراکنش استفاده کنید.

برای پاسخ به این سوال، ممکن است با جمع آوری اطلاعات در مورد هر خرید، مانند مبلغ، تاریخ، مکان، نوع خرید و آدرس دارنده کارت شروع کنید. شما به نمونه های زیادی از تراکنش ها، از جمله این اطلاعات، و همچنین برچسبی نیاز دارید که معتبر یا تقلبی بودن هر تراکنش را به شما بگوید . خوشبختانه، شما احتمالاً این اطلاعات را در یک پایگاه داده دارید. این رکوردها "داده های آموزشی ( training data )" نامیده می شوند و برای ساخت یک الگوریتم استفاده می شوند. هر بار که یک تراکنش جدید رخ می دهد، اطلاعات الگوریتم خود مانند مبلغ و تاریخ را می دهید و به سوال اصلی پاسخ می دهید: احتمال تقلبی بودن ( fraudulent ) این تراکنش چقدر است؟

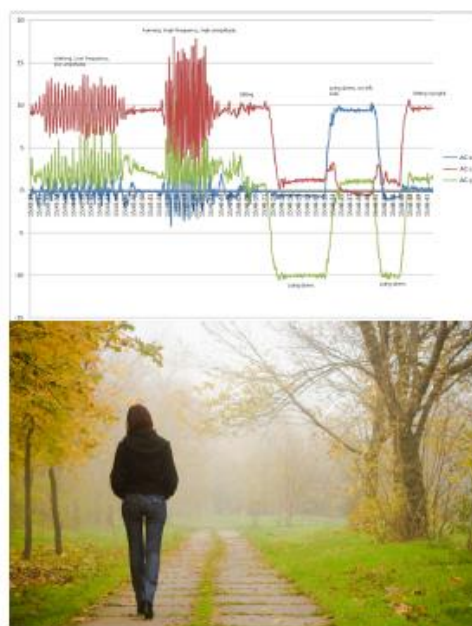


Amount	Date	Type	...
*	*	*	*
*	*	*	*
*	*	*	*
*	*	*	*
*	*	*	*
*	*	*	*
*	*	*	*

قبل از اینکه بتوانیم به این سوال پاسخ دهیم، بیایید مثال خود را مرور کنیم و آنچه را که برای یادگیری ماشینی نیاز داریم تا جادوی خود را انجام دهد برجسته کنیم. اول، یک مسئله علم داده با یک سوال کاملاً تعریف شده شروع می شود. سوال ما این بود که "احتمال تقلبی بودن این تراکنش چقدر است؟" در مرحله بعد، ما به برخی داده ها برای تجزیه و تحلیل نیاز داریم. ماه ها تراکنش های قدیمی کارت اعتباری و ابرداده های مرتبط داشتیم که قبلاً تقلبی یا معتبر شناخته شده بودند. در نهایت، هر بار که می خواهیم پیش بینی و قضاوت جدیدی انجام دهیم، به داده های اضافی نیاز داریم. ما باید در مورد هر خرید جدید اطلاعات یکسانی داشته باشیم تا بتوانیم آن را به عنوان "تقلبی" یا "معتبر" بدانیم.

## 2-2- اینترنت اشیا

حال، فرض کنید قصد دارید یک ساعت هوشمند برای نظارت بر فعالیت بدنی بسازید. شما می خواهید بتوانید فعالیت های مختلف مانند راه رفتن یا دویدن را به طور خودکار تشخیص دهید. ساعت هوشمند شما مجهز به حسگر ویژه ای به نام «شتاب سنج» است که حرکت را به صورت سه بعدی کنترل می کند. داده های تولید شده توسط این سنسور اساس یادگیری ماشین شما است. می توانید از چندین داوطلب بخواهید که ساعت شما را بپوشند و فعالیت خود را هنگام دویدن یا راه رفتن ضبط کنند. سپس می توانید الگوریتمی ایجاد کنید که داده های شتاب سنج را به عنوان نماینده یکی از این دو حالت تشخیص می دهد: راه رفتن یا دویدن.



ساعت هوشمند شما بخشی از یک زمینه در حال رشد به نام "اینترنت اشیا" است که به عنوان IOT نیز شناخته می شود، که اغلب با علم داده ترکیب می شود. اینترنت اشیا به ابزارهایی اطلاق می شود که کامپیوتر های استاندارد نیستند، اما همچنان توانایی انتقال داده ها را دارند. این شامل ساعت های هوشمند، سیستم های امنیتی خانگی متصل به اینترنت، سیستم های جمع آوری عوارض الکترونیکی، سیستم های مدیریت انرژی ساختمان و بسیاری موارد دیگر می شود. داده های تولید شده توسط اینترنت اشیا منبع عالی برای پروژه های علم داده است!

## 2-3- یادگیری عمیق

بیایید به مثال دیگری بپردازیم. یکی از وظایف کلیدی برای خودروهای خودران، شناسایی زمانی است که تصویری حاوی یک انسان است. مجموعه داده برای این مشکل چه خواهد بود؟



ما می توانیم تصویر را به صورت ماتریسی از اعداد بیان کنیم که در آن هر عدد نشان دهنده یک پیکسل است. با این حال، اگر ماتریس را به یک مدل یادگیری ماشین وارد کنیم، احتمالاً این رویکرد شکست خواهد خورد. به سادگی داده های ورودی بیش از حد وجود دارد!

1	1	1	1	1	1	1	1	2	2	3	3	3	1	3	1	1	2
1	1	1	1	1	1	1	1	2	2	3	3	3	1	3	1	1	2
1	1	1	1	1	1	1	1	2	1	2	3	1	3	3	1	1	2
8	8	8	8	8	8	8	8	8	8	8	8	48	20	20	20	8	8
6	6	7	6	6	6	6	6	6	6	6	6	6	20	20	20	8	5
4	4	4	4	4	4	4	4	4	4	4	4	4	20	20	20	5	5
4	4	4	4	4	4	4	5	5	5	5	5	5	5	4	4	4	4
4	4	4	4	5	5	5	5	5	5	4	4	4	4	4	4	4	4
5	5	5	5	5	5	4	4	4	4	4	4	4	4	4	4	4	4
4	4	5	5	5	4	4	4	4	4	4	4	4	4	4	4	4	4

ما به الگوریتم پیشرفته تری نیاز داریم که به عنوان یادگیری عمیق شناخته می شود. در یادگیری عمیق ، چندین لایه از مینی الگوریتم‌ها به نام « نورون‌ها » با هم کار می کنند تا نتیجه گیری های پیچیده ای را به دست آورند. یادگیری عمیق نسبت به یک مدل یادگیری ماشینی به داده های آموزشی بسیار بسیار بیشتری نیاز دارد، اما همچنین می تواند روابطی را بیاموزد که مدل های سنتی نمی توانند. یادگیری عمیق برای حل مشکلات پیچیده ای مانند طبقه بندی تصویر یا درک زبان استفاده می شود.

### 3- ساخت تیم علوم داده

در ادامه شما یاد خواهید گرفت که چگونه تیم داده خود را برای برآوردن نیازهای سازمان خود بسازید و ساختار دهید.

#### 3-1-1- اعضای تیم

شاید تعجب کنید که بدانید "علم داده" یک رشته واحد نیست. این در واقع سه شغل مختلف است: مهندس داده، تحلیلگر داده، و دانشمند یادگیری ماشین. بیایید هر یک را بررسی کنیم.



*Data Engineer*



*Data Analyst*



*Machine Learning  
Scientist*

#### 3-1-1-1- مهندس داده

مهندسان داده جریان اطلاعات را کنترل می کنند: آنها سیستم های ذخیره سازی داده ها و زیرساخت های تخصصی ایجاد می کنند تا اطمینان حاصل کنند که داده ها به راحتی به دست می آیند و پردازش می شوند.



*Data Engineer*

- Information architects
- Build storage solutions
- Maintain data access

اکثر مهندسان داده با SQL آشنا هستند که از آن برای ذخیره و مدیریت کلان داده ها استفاده می کنند. آنها همچنین از یکی از زبان های برنامه نویسی مانند جاوا، اسکالا یا پایتون برای پردازش داده ها و خودکارسازی وظایف مربوطه استفاده می کنند.



## Data Engineer

- SQL
  - Storing large quantities of data
- Java, Scala, or Python
  - Programming languages for processing data and automating tasks

### 3-1-2- تحلیلگر داده

تحلیلگران داده ، حال را از طریق داده ها توصیف می کنند. آنها این کار را با داشبورد، آزمون فرضیه ها ( hypothesis tests ) و تجسم ( visualization ) انجام می دهند. آنها اغلب پیشینه ای در زمینه آمار یا علوم کامپیوتر دارند، اما تمایل دارند تجربه مهندسی کمتری نسبت به مهندسان داده و تجربه ریاضی کمتری نسبت به دانشمندان داده داشته باشند.



## Data Analyst

- Creating dashboards
- Hypothesis testing
- Data visualization

تحلیلگران داده از spreadsheets برای انجام تجزیه و تحلیل های ساده بر روی مقادیر کم داده استفاده می کنند. آنها از SQL ، همان زبانی که توسط مهندسان داده استفاده می شود، برای تحلیل های بزرگتر استفاده می کنند. در حالی که مهندسان داده راه حل های ذخیره سازی SQL را می سازند و پیکربندی می کنند، تحلیلگران داده از پایگاه های داده موجود برای مصرف و خلاصه کردن داده ها استفاده می کنند. تحلیلگران همچنین از هوش تجاری یا BI ، ابزارهایی مانند Tableau ، Power BI یا Looker برای ایجاد داشبورد و به اشتراک گذاری تحلیل های خود استفاده می کنند.





## Data Analyst

- Spreadsheets (Excel or Google Sheets)
  - Simple storage and analysis
- SQL
  - Large-scale analysis
- BI Tools (Tableau, Power BI, Looker)
  - Dashboarding and sharing information

### 3-1-3 — دانشمند یادگیری ماشین

یادگیری ماشین شاید پرهیاهوترین بخش علم داده باشد. از آن برای تشخیص و حدس آنچه که احتمالاً درست است از آنچه قبلاً می دانیم استفاده می شود. این دانشمندان از داده های آموزشی برای طبقه بندی داده های بزرگتر و غیر طبقه بندی شده استفاده می کنند. یادگیری ماشینی می تواند به ما بگوید ارزش یک سهام در هفته آینده چقدر است، کدام تصاویر حاوی یک ماشین هستند یا چه احساساتی توسط مجموعه ای از توییت ها بیان می شوند.



## Machine Learning Scientist

- Predictions and extrapolations
- Classification
- Stock price prediction
- Image processing
- Automated text analysis

دانشمندان یادگیری ماشین از پایتون یا R برای ایجاد مدل های پیش بینی خود استفاده می کنند. هر دو زبان برنامه نویسی عالی برای علم داده هستند و فردی که یک زبان را می داند احتمالاً می تواند کد را در زبان دیگر بخواند. به یاد داشته باشید، یادگیری زبان های برنامه نویسی به اندازه زبان های گفتاری دشوار نیست. اگر کسی بلد باشد فرانسوی صحبت کند، ممکن است سالها طول بکشد تا زبان آلمانی را یاد بگیرد. زبان های برنامه نویسی بیشتر شبیه ابزارها هستند. اگر می دانید چگونه از دریل برقی استفاده کنید، لزوماً نحوه استفاده از اره برقی را بلد نیستید، اما احتمالاً می توانید با کمی آموزش یاد بگیرید!



## Machine Learning Scientist

- Python and R
  - Programming languages for creating predictive models

### 3-1-4- مرور

برای جمع‌بندی: مهندسان داده داده‌ها را ذخیره و نگهداری می‌کنند، تحلیلگران داده‌ها داده‌ها را تجسم (visualize) و توصیف می‌کنند، و دانشمندان یادگیری ماشین با داده‌ها مدل‌سازی و پیش‌بینی می‌کنند. هر موقعیت از مجموعه کمی متفاوت از ابزارها برای رسیدن به اهداف خود استفاده می‌کند.

Data Engineer	Data Analyst	Machine Learning Scientist
Store and maintain data	Visualize and describe data	Model and predict with data
SQL + Java/Scala/Python	SQL + BI Tools + Spreadsheets	Python/R

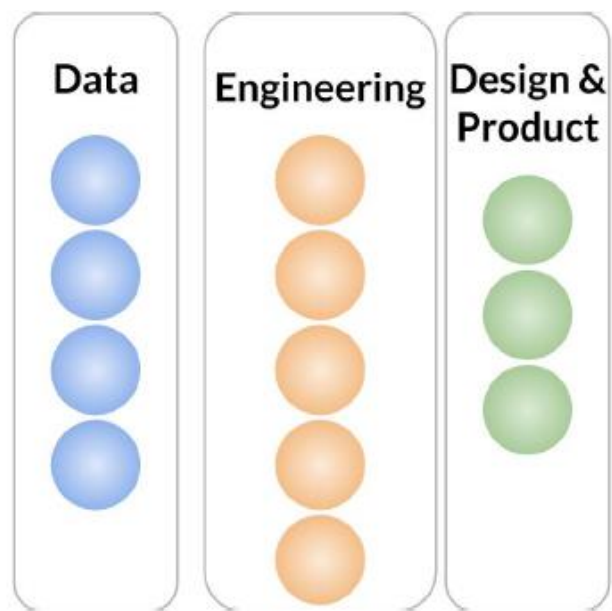
### 3-2- ساختار تیم علم داده

هنگامی که تعدادی متخصص داده را استخدام کردید، سه راه اصلی وجود دارد که می‌توانید ساختار تیم داده خود را سامان دهید: ایزوله شده ، تعبیه شده یا ترکیبی.

#### 3-2-1- ایزوله شده (isolated)

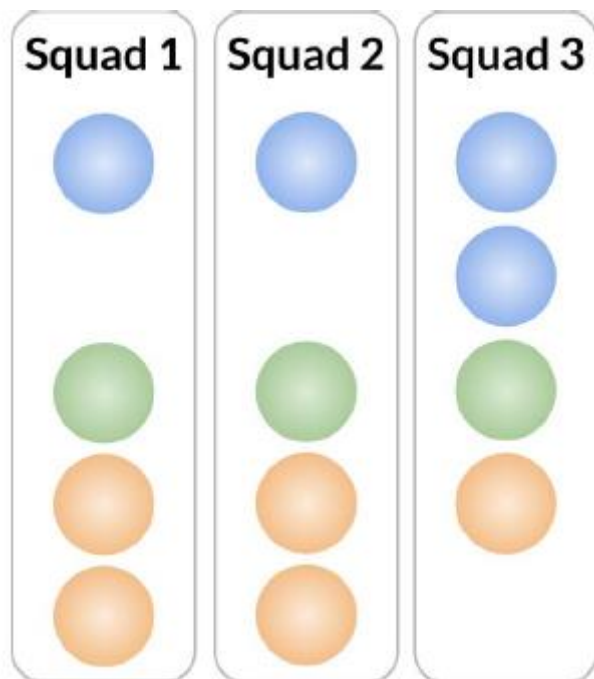
یک تیم داده ایزوله می‌تواند شامل یک یا چند نوع کارمند داده بدون هیچ تیم دیگری مانند مهندسی یا محصول باشد. این یک ساختار عالی برای آموزش اعضای جدید تیم و تغییر سریع پروژه ای است که هر عضو روی آن کار می‌کند.





### 3-2-2- تعبیه شده (embedded)

روش دیگر، استفاده از یک مدل تعبیه شده است که در آن هر کارمند داده بخشی از یک تیم است، که شامل مهندسان و مدیران محصول نیز می شود، این روش می تواند مفید باشد. این مدل به هر کارمند داده ای اجازه می دهد تا در یک پروژه تجاری خاص تجربه کسب کند و تبدیل به یک متخصص ارزشمند شود.

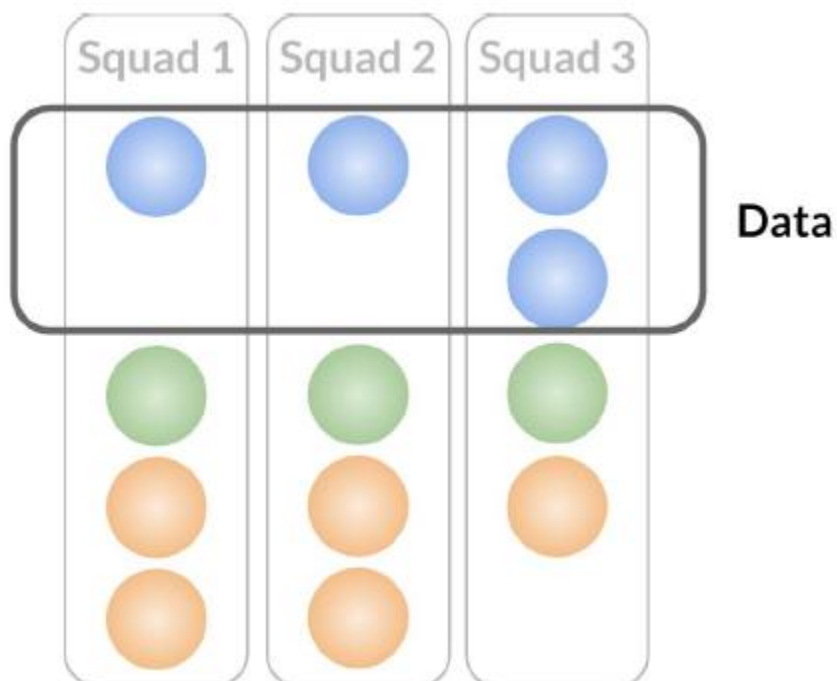


### 3-2-3 ترکیبی (hybrid)

مدل هیبریدی شبیه به مدل تعبیه شده به نظر می‌رسد، اما با یک همگام‌سازی اضافی برای همه کارکنان داده در تمام

گروه‌ها. این لایه اضافی سازمان، بدون در نظر گرفتن اینکه کارمند به کدام پروژه منصوب شده است، امکان فرآیندهای داده

یکنواخت و توسعه شغلی را فراهم می‌کند.



## 4- منابع داده و خطرات

قبلاً با نحوه کار در علم داده آشنا شدید. در این درس، ما بر مرحله اول تمرکز خواهیم کرد: جمع آوری داده ها.

### 4-1- منابع مشترک داده ها

داده ها در همه جا وجود دارد و تقریباً هر فرآیند تجاری می تواند کوه هایی از داده را تولید کند. برخی از رایج ترین منابع داده، رویدادهای وب، داده های مشتری، داده های لجستیک و تراکنش های مالی هستند. این امکان وجود دارد که شرکت شما در حال حاضر تمام این اطلاعات را جمع آوری کرده باشد. بهتر است از مهندسان داده خود بپرسید که چه چیزی جمع آوری می شود و چه چیزی جمع آوری نمی شود و بر اهمیت شروع فرآیند جمع آوری داده های متمرکز زودتر از دیرتر تاکید کنید.

## Common sources of data

- Web events
- Customer data
- Logistics data
- Financial transactions

### 4-2- داده های وب

بیایید نگاهی عمیق تر به داده های وب بکنیم. هنگامی که کاربر از یک صفحه وب بازدید می کند یا روی یک پیوند کلیک می کند، ردیابی این اطلاعات به منظور محاسبه نرخ تبدیل ( conversion rates ) یا نظارت بر محبوبیت قسمت های مختلف محتوا می تواند مفید باشد. حداقل، می خواهید نام رویداد را جمع آوری کنید، که می تواند به معنای نشانی اینترنتی صفحه بازدید شده یا شناسه عنصری که روی آن کلیک شده، مهر زمانی رویداد ( timestamp )، و یک شناسه برای کاربری باشد که این کار را انجام داده است.

## Web data

- Events
- Timestamps
- User information

user_id	event_name	timestamp
1234	homepage_visit	2019-01-01 12:01:01

### 4-3- اطلاعات قابل شناسایی شخصی (PII)

فرض کنید علی محمدی مشتری است که از وب سایت شرکت شما بازدید می کند و یکی از محصولات شما را می پسندد. ممکن است انتخاب کنید نام او، مهر زمانی و شیئی که روی آن کلیک کرده است را ردیابی کنید. مهم است که به یاد داشته باشید که نام علی محمدی، اطلاعات قابل شناسایی شخصی (Personally Identifiable Information) یا PII است. PII شامل نام، موقعیت مکانی، آدرس ایمیل، و هر اطلاعات دیگری است که می تواند برای پیوند دادن یک رویداد وب به یک انسان واقعی استفاده شود. بایستی با PII با احتیاط زیادی برخورد کرد.

Timestamp	Object Clicked
2019-01-20 12:05:00	Like Button

### 4-4- نام مستعار داده ها (Data pseudonymization)

یکی از ساده ترین راهها برای محافظت از هویت علی محمدی، تقسیم این اطلاعات به دو ورودی جداگانه است. ما می توانیم به علی یک شناسه کاربری، در این مورد 185477 اختصاص دهیم و آن اطلاعات را در جدول کاربران ذخیره کنیم. سپس می توانیم رویداد او را با استفاده از این شناسه شناسایی کنیم. ما داده های جدول رویدادها را با نام مستعار می خوانیم زیرا علی را نمی توان به تنهایی با آن جدول شناسایی کرد، اما اگر اطلاعات جدول کاربران را با جدول رویدادها ترکیب کنیم، می توان او را شناسایی کرد. برای محافظت از علی، می خواهیم مطمئن شویم که دسترسی به جدول کاربران فقط به افرادی محدود می شود که

باید هویت علی را بدانند، مانند نمایندگان ارشد تیم خدمات مشتری یا اعضای تیم حقوقی. همچنین می‌خواهیم به‌طور دوره‌ای بازرسی کنیم که چه کسی به این داده‌ها دسترسی داشته است و چگونه از آن استفاده کرده است تا اطمینان حاصل شود که داده‌های علی امن هستند و به حریم شخصی او احترام گذاشته می‌شود.

user_id	Timestamp	Object Clicked
185477	2019-01-20 12:05:00	Like Button

#### 4-5- ناشناس سازی داده‌ها (Data anonymization)

بهترین راه برای محافظت از حریم خصوصی علی این است که پس از اختصاص شناسه کاربری علی، اطلاعات جدول کاربران را از بین ببرید. بدون جدول کاربران، جدول رویدادها داده‌های کاملاً ناشناس است. برای بسیاری از اهداف تحلیل، داده‌های ناشناس کافی است. ما باید بدانیم که جین یک فرد منحصر به فرد است، اما نیازی نیست که نام او یا هر PII دیگری را بدانیم.

#### 4-6- مقررات عمومی حفاظت از داده‌ها

ممکن است اخیراً اصطلاح "GDPR" را از تیم داده خود شنیده باشید. GDPR مخفف عبارت General Data Protection Regulation است و برای تمام داده‌های داخل اتحادیه اروپا اعمال می‌شود. هدف از GDPR این است که افراد بر داده‌های شخصی خود کنترل داشته باشند. از جمله موارد دیگر، GDPR مدت زمان ذخیره داده‌ها را تنظیم می‌کند، ناشناس سازی مناسب را الزامی می‌کند و نیاز به افشای جمع‌آوری داده‌ها و کسب رضایت دارد.

## General Data Protection Regulation (GDPR)

- Applies to all data inside of the EU
- Give individuals control over their personal data
- Regulates how long data can be stored
- Mandates appropriate anonymization
- Disclose data collection and gain consent

## 5- داده های درخواستی ( Solicited data )

در چند صفحه گذشته، در مورد داده هایی که شرکت شما می تواند از طریق وب یا تراکنش های مالی جمع آوری کند، یاد گرفتیم. اکنون، ما داده هایی را پوشش می دهیم که می توانید با پرسیدن نظرات مشتریان خود به دست آورید، که به آن داده های درخواستی می گوئیم.

### 5-1- چرا داده ها را درخواست می کنیم؟

از داده های درخواستی می توان برای ایجاد وثیقه بازاریابی استفاده کرد، مانند پستی در مورد اینکه چند درصد از مشتریان ما از محصول ما راضی هستند. همچنین می توان از آن برای کاهش ریسک فرآیند تصمیم گیری استفاده کرد، مانند زمانی که از کاربران نظرسنجی می کنیم تا علاقه به محصول جدید را ارزیابی کنیم. در نهایت، می توان از آن برای نظارت بر کیفیت استفاده کرد.



### 5-2- انواع داده ای درخواستی؟

انواع متداول داده های درخواستی شامل نظرسنجی ها، بررسی های مشتریان، پرسشنامه های درون برنامه ای و گروه های تمرکز است. این تصویر نوع بسیار رایج بازخورد درخواستی را نشان می دهد: امتیاز خالص تبلیغ کننده یا NPS، که از کاربر می پرسد چقدر احتمال دارد که کاربر یک محصول را به یک دوست یا همکار توصیه کند.

داده های درخواستی می تواند کیفی باشد، مانند مکالمات و سؤالات باز، یا کمی باشد، مانند سؤالات چند گزینه ای یا مقیاس های رتبه بندی. داده های کیفی بسیار ذهنی هستند و نیاز به تحلیل زیادی دارند. داده های کمی را می توان به راحتی در یک نمودار یا نمودار خلاصه کرد. به طور کلی، جمع آوری داده های کیفی در مقیاس کوچک برای ایجاد فرضیه خوب است. به عنوان مثال، یک گروه تمرکز ممکن است ایده هایی در مورد ویژگی هایی که ممکن است بخواهیم بسازیم ارائه دهد. برای تایید این فرضیه ها به مجموعه کمی در مقیاس بزرگتر نیاز است. به عنوان مثال، می توانیم از کاربران بخواهیم فهرستی از ویژگی ها را از مطلوب ترین به حداقل مطلوب رتبه بندی کنند.

### 5-3- ترجیحات آشکار و بیان شده (Revealed and stated preferences)

مهم است که به یاد داشته باشید که داده های درخواست شده به طور کلی اولویت اعلام شده کاربران را به ما می گوید. ترجیح اعلام شده چیزی است که کسی به ما می گوید که می خواهد یا معتقد است و تا حدودی فرضی است. هنگامی که کاربر واقعاً اقدامی مانند خرید یک محصول انجام می دهد، ترجیحات آشکار او را می آموزیم. ما امیدواریم که اولویت های اعلام شده کاربران ما، شاخص های خوبی برای اولویت آشکار آنها باشد، اما همیشه اینطور نیست. بسیاری از مردم ترجیح می دهند مرتباً به باشگاه بروند. با این حال، ترجیح بسیاری از همان افراد نشان داده شده تنها رفتن گهگاهی است. برخی از باشگاه ها یک مدل کسب و کار کامل دارند که بر اساس تفاوت مورد انتظار بین اولویت های اعلام شده و آشکار افراد برای ورزش است.

اکنون که انواع داده‌های درخواستی را می‌شناسیم و از مشکلات اولویت‌های آشکار شده در مقابل اولویت‌های اعلام‌شده آگاه هستیم، اجازه دهید برخی از بهترین شیوه‌ها را مرور کنیم. اول، ما باید سعی کنیم در هنگام پرسیدن سوال تا حد امکان دقیق باشیم. این ویژگی باید هم در مورد جمله بندی سؤال و هم در مورد پاسخ های بالقوه ای که ارائه می دهیم اعمال شود.

#### 5-4- بهترین شیوه ها

در مرحله بعد، ما باید از زبان بارگذاری شده ( loaded language ) اجتناب کنیم، به خصوص اگر ممکن است پاسخ دهندگان را نسبت به یک انتخاب خاص سوگیری کند. به عنوان مثال، از استفاده از صفت برای توصیف گزینه های احتمالی خودداری کنید و سعی کنید عینی ( objective ) باشید.

در صورت امکان، نظرسنجی خود را با مقایسه با مقادیر شناخته شده کالیبره کنید. به عنوان مثال، به جای اینکه از پاسخ دهنده بپرسید که آیا به یک محصول جدید علاقه مند است، از او بخواهید که علاقه خود را به آن محصول با علاقه خود به یکی از پیشنهادات فعلی شما یا پیشنهادات یک رقیب شناخته شده مقایسه کند.

در نهایت، پرسیدن هر چه بیشتر سؤالات ممکن است وسوسه انگیز باشد. با این گزینه مبارزه کنید و در عوض اطمینان حاصل کنید که هر سؤالی که می‌پرسید به شما کمک می‌کند تا یک اقدام قاطع انجام دهید. به عنوان مثال، اگر مشتریان ترجیحی برای ویژگی A نسبت به ویژگی B نشان دادند، شما باید ابتدا متعهد شوید که ویژگی A را بسازید.



## 6- جمع آوری داده های اضافه ( additional data )

در حالی که جمع آوری داده های داخلی برای برخی از پروژه های علم داده مفید است، اما تنها یک تکه از این پازل است.

اغلب، شما نیاز به جمع آوری داده ها از منابع خارجی نیز دارید.



### 6-1- حتی داده های بیشتر

راه های زیادی وجود دارد که می توانید داده های اضافی را برای سازمان خود جمع آوری کنید. چند راه متداول شامل

API ها، رکوردهای عمومی و Mechanical Turk است که در ادامه به همه آنها خواهیم پرداخت.

#### 6-1-1- API های داده

بیایید با API ها شروع کنیم. API مخفف Application Programming Interface است. این یک راه آسان برای

درخواست داده از شخص ثالث از طریق اینترنت است. بسیاری از شرکت ها API دارند تا به تیم شما اجازه دسترسی به داده هایشان

را بدهند. برخی از API های قابل توجه عبارتند از: Wikipedia، Yahoo!، و نقشه های گوگل، اما بسیاری از موارد

دیگر وجود دارد. اگر با یک شریک کار می کنید و فکر می کنید که آنها ممکن است داده های مفیدی داشته باشند، یک جستجوی

سریع در وب انجام دهید و ببینید آیا یک API وجود دارد یا خیر!



بیایید به نمونه ای از API توییتر نگاه کنیم. فرض کنید می‌خواهیم توییت‌ها را با `#Example` دنبال کنیم. ما می‌توانیم از API توییتر برای درخواست همه توییت‌های دارای این هشتگ استفاده کنیم. در این مرحله، ما گزینه‌های زیادی برای تجزیه و تحلیل داریم. ما می‌توانیم روی متن هر توییت یک تحلیل احساسی انجام دهیم و درباره نظرات و احساسات مردم ایده بگیریم. ما به سادگی می‌توانیم ردیابی کنیم که هر هفته چند بار هشتگ `Example` ظاهر می‌شود.

#### 6-1-2- سوابق عمومی

سوابق عمومی یکی دیگر از راه‌های عالی برای جمع‌آوری داده‌های اضافی است. در ایالات متحده، `data.gov` داده‌های بهداشت، آموزش و تجارت را برای دانلود رایگان در دسترس دارد. در اتحادیه اروپا، `data.europa.eu` داده‌های مشابهی دارد. اینها می‌توانند منابع خوبی برای درک روندهای سطح جمعیت (population-level trends) یا جمع‌آوری داده‌های اقتصادی باشند.



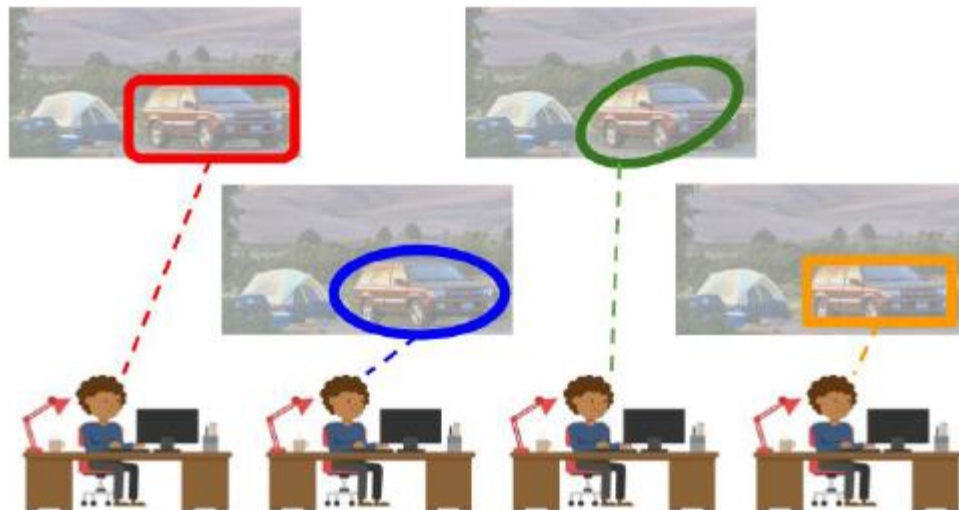
### 6-1-3- ترک مکانیکی

قبلاً، ما در مورد تشخیص تصویر به عنوان یک نوع مشکل علم داده بحث کردیم. برای ساختن یک الگوریتم تشخیص تصویر خوب، به مجموعه ای از تصاویر نیاز داریم که در آن تصاویر قبلاً برچسب گذاری شده باشند که به آن مجموعه آموزشی (training set) می گویند. اما ما فقط به یک یا دو عکس نیاز نداریم. ما به صدها یا هزاران عکس نیاز داریم. دریافت این تصاویر برچسب گذاری شده می تواند واقعاً سخت و زمان بر باشد و فقدان مجموعه آموزشی اغلب باعث می شود پروژه های علمی داده خوب تکمیل نشوند.

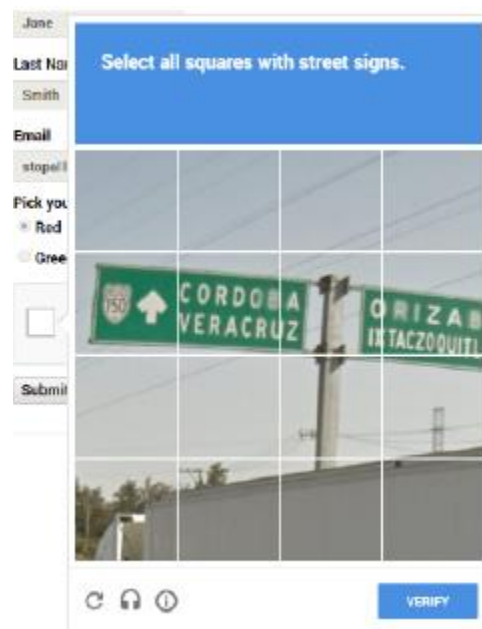


بسته به نوع مجموعه آموزشی مورد نیاز، Mechanical Turk که MTurk نیز نامیده می شود، می تواند یک گزینه عالی باشد. ترک مکانیکی به معنای درخواست از انسان برای تکمیل کاری است که ما در نهایت برای کامپیوتری کردن آن عمل برنامه ریزی می کنیم. در مثال قبلی ما، این به معنای برچسب زدن تعدادی عکس برای ایجاد یک مجموعه آموزشی برای تشخیص تصویر است. به جای اینکه از یک نفر بخواهیم هزاران تصویر را برچسب گذاری کند، هزاران نفر را استخدام می کنیم و به هر یک از آنها برای برچسب زدن چند تصویر پول می دهیم. برای اطمینان از کیفیت، ممکن است از دو یا سه نفر بخواهیم تصویر مشابه را مرور کنند و سپس رایج ترین پاسخ را در نظر می گیریم .

# Select the car in the image.



پلتفرم‌های زیادی برای کمک به Mechanical Turk و استخدام کمک‌کنندگان مانند AWS MTurk وجود دارد. ترک مکانیکی فقط برای تشخیص تصویر نیست. همچنین می‌توانید از آن برای برجسب زدن نظرات مشتریان به عنوان مثبت یا منفی، استخراج متن از فرم یا برجسته کردن کلمات کلیدی در یک جمله استفاده کنید. در مثال سمت راست، از کاربران خواسته می‌شود که مشخص کنند کدام بخش از تصویر دارای علامت خیابان است.



## 7- ذخیره و بازیابی داده ها

تا اینجا، با بسیاری از منابع داده آشنا شدید. اکنون، بیایید روش های کارآمد ذخیره سازی داده هایی را که سازمان شما ممکن است جمع آوری کند، مورد بحث قرار دهیم.

### 7-1- راه حل های ذخیره سازی موازی

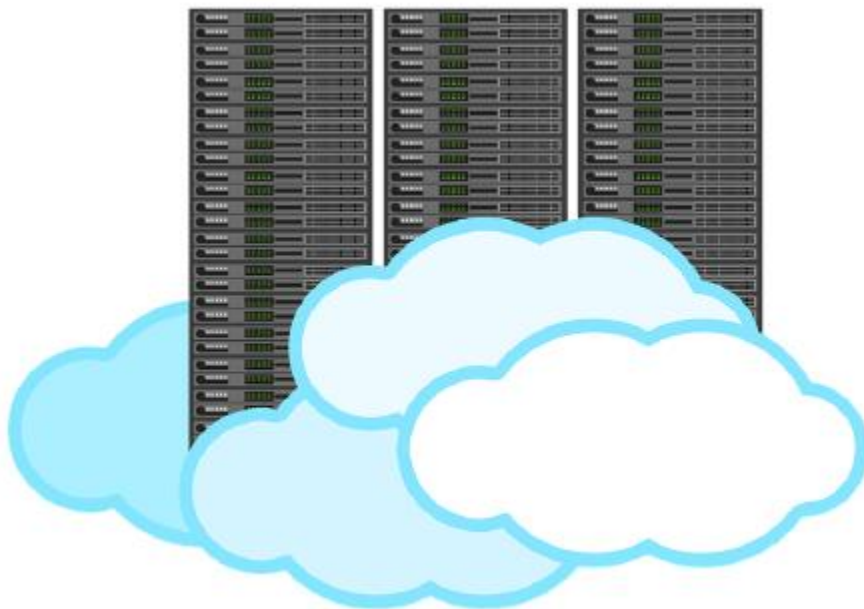
کسب و کار شما احتمالاً داده های بسیار بیشتری از آنچه که حتی در یک رایانه ذخیره می شود تولید می کند. برای اطمینان از اینکه همه داده ها ذخیره شده و دسترسی به آنها آسان است، باید آن ها را در رایانه های مختلف ذخیره کنید. شرکت شما ممکن است مجموعه ای از رایانه های ذخیره سازی خود را به نام "خوشه" یا "سرور" در محل داشته باشد.



### 7-2- ابر

از طرف دیگر، ممکن است شرکت شما به شرکت دیگری پول بدهد تا داده ها را برای شما ذخیره کند. این به عنوان

"ذخیره سازی ابری" نامیده می شود. ارائه دهندگان فضای ابری رایج عبارتند از Amazon Web Services، Microsoft Azure و Google Cloud. این خدمات بیش از ذخیره سازی داده ها را ارائه می دهند. آنها همچنین می توانند به سازمان شما در تجزیه و تحلیل داده ها، یادگیری ماشینی و یادگیری عمیق کمک کنند. در حال حاضر، ما فقط بر روی ذخیره سازی داده ها تمرکز می کنیم.



### 7-3- انواع ذخیره سازی داده ها

انواع مختلف داده ها به راه حل های ذخیره سازی متفاوتی نیاز دارند. برخی از داده ها مانند ایمیل، متن، فایل های ویدیویی و صوتی، صفحات وب و پیام های رسانه های اجتماعی بدون ساختار هستند. این نوع داده ها اغلب در نوعی پایگاه داده به نام پایگاه داده اسناد ذخیره می شوند.

#### Unstructured

- Email
- Text
- Video and audio files
- Web pages
- Social media

#### Document Database

معمولاً داده ها را می توان به صورت جداول اطلاعات بیان کرد، مانند آنچه ممکن است در یک صفحه گسترده بیابید.

پایگاه داده ای که اطلاعات را در جداول ذخیره می کند، پایگاه داده رابطه ای نامیده می شود. هر دوی این نوع پایگاه های داده را می توان در ارائه دهندگان ذخیره سازی ابری که قبلاً ذکر شد، یافت.

## Tabular

Customer Name	Customer Address	...
Jane Doe	123 Maple St.	...

## Relational Database

### 7-4- پرس و جو داده ها

هنگامی که داده ها در یک پایگاه داده سند یا یک پایگاه داده رابطه ای ذخیره شدند، باید به آن دسترسی داشته باشیم.

در سطح پایه، ما می خواهیم بتوانیم یک داده خاص مانند «همه تصاویری که در 3 بهمن ماه ایجاد شده اند» یا «همه آدرس های مشتریان در شهر تهران» را درخواست کنیم. علاوه بر این، حتی ممکن است بخواهیم تحلیل هایی مانند جمع، شمارش یا میانگین گیری داده ها را انجام دهیم.





هر نوع پایگاه داده زبان پرس و جو خود را دارد. پایگاه‌های داده اسناد عمدتاً از NoSQL استفاده می‌کنند، در حالی که پایگاه‌های داده رابطه‌ای عمدتاً از SQL استفاده می‌کنند.

Data Type	Query Language
Document Database	NoSQL
Relational Database	SQL

### 7-5- کنار هم قرار دادن همه : مکان

ذخیره داده‌های شرکت شما مانند ساخت یک کتابخانه است. ابتدا باید تصمیم بگیرید که کجا کتابخانه خود را بسازید. این مربوط به انتخاب یک ابر است: یا یک خوشه داخلی یا یکی از ارائه‌دهندگانی که قبلاً در مورد آن صحبت کردیم، Azure ، AWS یا Google Cloud .



### 7-6- کنار هم قرار دادن همه : نوع داده‌ها

در مرحله بعد، باید تصمیم بگیرید که چه نوع قفسه‌هایی را برای نگهداری کتاب‌های خود نصب کنید. انواع قفسه‌ها به انواع کتاب‌ها بستگی دارد.





این مشابه انتخاب بین یک پایگاه داده سند برای داده های بدون ساختار یا یک پایگاه داده رابطه ای برای داده های جدولی است. درست مانند یک کتابخانه که ممکن است چندین نوع قفسه داشته باشد، ممکن است نیاز باشد که برخی از داده ها در یک پایگاه داده اسناد و سایر داده ها در یک پایگاه داده رابطه ای ذخیره شوند.

Data Type	Storage Solution
Unstructured	Document Database
Tabular	Relational Database

## 7-7- کنار هم قرار دادن همه : پرس و جو

در نهایت، به سیستمی برای ارجاع و بررسی کتاب ها نیاز دارید. نحوه مکان یابی و بازیابی هر کتاب به نحوه ذخیره آن کتاب بستگی دارد.



به طور مشابه، برای صحبت با پایگاه داده به یک زبان پرس و جو نیاز دارید. برای پایگاه‌های داده اسناد، ما معمولاً از NoSQL و برای پایگاه‌های داده رابطه‌ای، معمولاً از SQL استفاده می‌کنیم.

Data Type	Query Language
Document Database	NoSQL
Relational Database	SQL

## منابع

- Data Science for Business and Decision Making by Luiz Favero and Patrícia Belfiore .1
- Harvard Business School at online.hbs.edu .2
- Harvard Business Review at hbr.org .3
- wikipedia.org .4
- ibm.com/cloud/learn/data-science-introduction .5
- coursera.org/browse/data-science .6
- <http://cloudscaling.com/blog/cloud-computing> .7
- <https://www.redhat.com/en/technologies/management> .8
- <https://github.com> .9
- <https://www.darkreading.com/cloud> .10