



Data Science



علم داده

فرصت ایجاد یک
کسب و کار جدید

ارائه:

حسین حق‌بین

مرداد ماه ۱۳۹۹

بِسْمِ اللَّهِ الرَّحْمَنِ الرَّحِيمِ

✓ علم داده‌ها چیست؟

✓ دانشمند داده به چه کسی گفته می‌شود؟

✓ موقعیت‌های شغلی در حوزه علم داده

✓ علم داده در استارت‌آپ‌ها

✓ استارت‌آپ‌های موفق در حوزه علم داده



علم داده‌ها چیست

علم داده‌ها (**Data Science**)، یک زمینه میان‌رشته‌ای است که از روش‌ها، فرآیندها، الگوریتم‌ها و سیستم‌های علمی برای استخراج دانش (**Knowledge**) از داده‌های خام (**Raw Data**) استفاده می‌کند.

در علم داده نظریه‌های علوم گوناگون از جمله
آمار، علوم کامپیوتر، ریاضیات، و دیگر روش‌ها
مانند یادگیری ماشین، داده‌کاوی و بصری‌سازی
داده‌ها مورد استفاده قرار می‌گیرد.

علم داده را پارادایم چهارم علم می‌دانند.

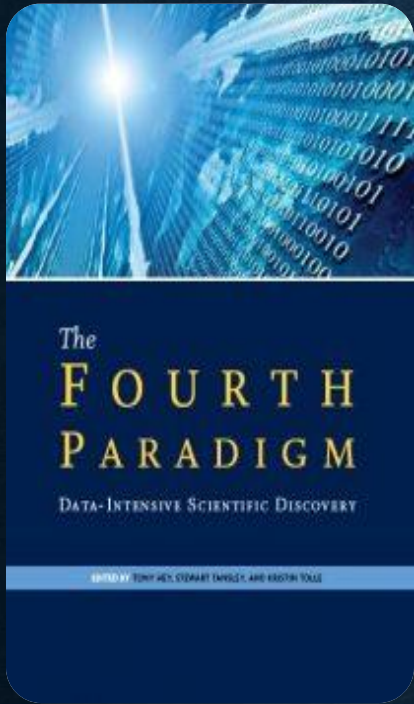
Tansley and Tolle (2009)

1. شواهد تجربی (Empirical Evidence)

2. نظریه علمی (Scientific Theory)

3. علوم محاسباتی (Computational Science)

4. علم داده (Data Science)



دانشمندان بر این باور هستند که

«کلیه موارد مربوط به علم تحت تاثیر

فناوری های **داده محور** در حال تغییر است».

مجمع جهانی اقتصاد در سال ۲۰۱۱ اعلام کرد که: ((داده‌ها نفت عصر جدید هستند))

“DATA IS THE NEW OIL.”

Coined in 2006 by Clive Humby, a British data commercialization entrepreneur, this now famous phrase was embraced by the World Economic Forum in a 2011 report, which considered data to be an economic asset, like oil.

From the beginning of recorded time until 2003, we created **5 exabytes** of data.

In 2011 the same amount was created every two days.

By 2013, it's expected that the time will shrink to 10 minutes.

7 billion DVDs.

Side by side, that's that's seven times the height of Everest.

There are nearly as many bits of information in the digital universe as there are stars in our actual universe.

As of August 2012, there were just over **4 million** articles in the English Wikipedia.

There are **133 million** BLOGS on the web.

80% of all humans own a mobile phone of some sort. Out of 5 billion mobiles, 1 billion are smartphones. (In Singapore, 84% of citizens are smartphone users.)

English is the dominant language of the web. But by 2014 it will be **Chinese**, if its current rate of increase continues.

Top languages used on the web (May 2011):



247 billion EMAILS are sent every day. (Up to 90% are spam.)

10% of all photos ever taken were taken in 2011.

60% of all humans (5.4 billion people) are active texters. In 2010, 153,000 text messages were sent every second.

50% of 8-year-old kids in the U.S. are given access to a smartphone.

Just as a study of activity on Twitter gave residents, family members, and journalists advance warning of details about the devastating earthquake and tsunami in Japan, **high-frequency traders**, with the help of computer algorithms, use Big Data to follow trends and to act quickly on their findings.

These specialized algorithms make split-second decisions to buy or sell a commodity. New cable being laid under the Atlantic will shave **5 milliseconds** from the current 85 milliseconds it takes for trading instructions to travel between New York City and London.

With new fiber-optic cable, the round-trip time between New York and London will be 80.6 milliseconds.

This 5-millisecond saving is worth many millions of dollars to the trading firms who use the cable (and who will pay millions to do so).

How they save 5 milliseconds

The depth of the Atlantic Ocean varies. The new cable will lie on areas of the ocean floor that are up to 1,000 feet shallower than the current fastest cable. By taking a different route, the new cable is shorter, meaning that the time it takes for messages to travel along it is shortened.



آینده متعلق افراد یا شرکتهایی است که
بتوانند داده‌ها را به محصول تبدیل کنند.

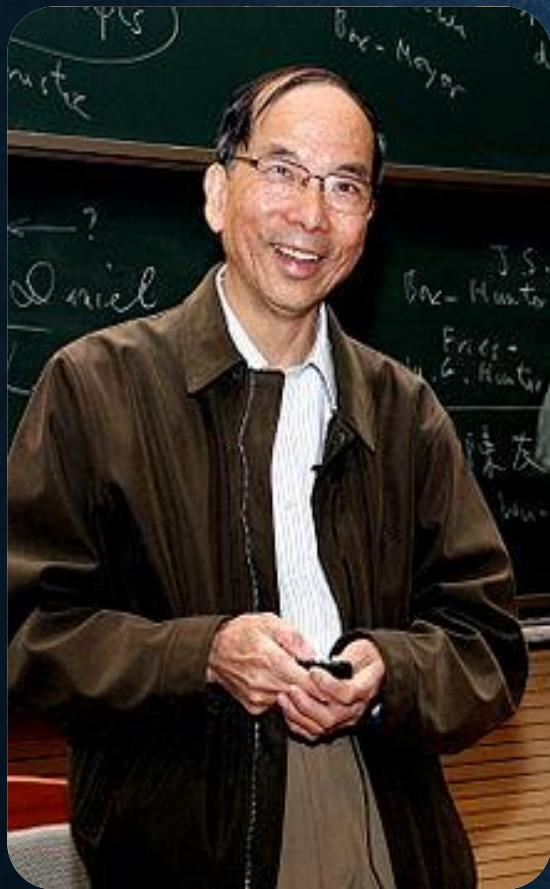
تاریخچه علم داده



در سال ۱۹۶۰، «پیتر ناور» (Peter Naur) از این عبارت به عنوان جایگزینی برای علوم کامپیوتر استفاده کرد. ناور بعدها اصطلاح «داده‌شناسی» (datalogy) را بدین منظور معرفی کرد. وی در سال ۱۹۷۴ در مقاله‌ای با عنوان «بررسی دقیق روش‌های کامپیوتری» از اصطلاح علم داده برای بیان پردازش‌های داده آن دوران استفاده کرد.

در سال ۱۹۹۶، اعضای «فدراسیون بین‌المللی جامعه دسته‌بندی» (IFCS) برای گردهمایی دو سال یکبار خود، در شهر کوبه ژاپن گردهم آمدند. در گردهمایی مذکور، برای اولین بار از اصطلاح علم داده به عنوان اسم کنفرانس علم داده، دسته‌بندی و روش‌های مرتبط – استفاده شد.

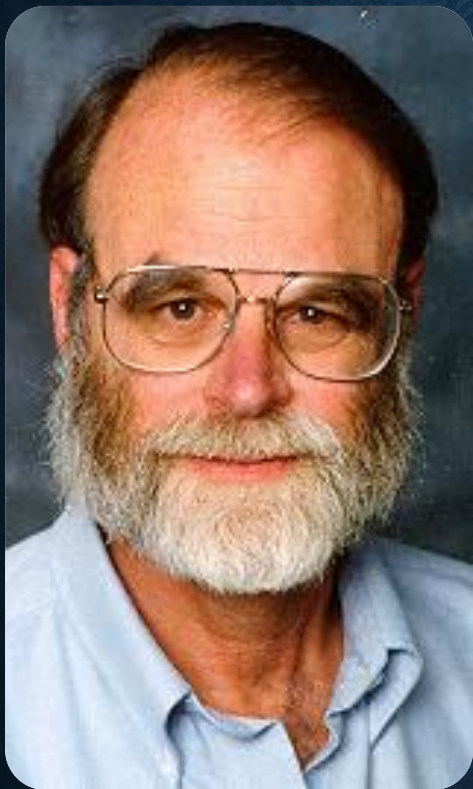




در نوامبر سال ۱۹۹۷، سی اف جف وو (C. F. Jeff Wu) ، آماردان و استاد تحلیلگر شرکت کوکا-کولا، سخنرانی افتتاحیه رویدادی در حوزه علم داده را با عنوان ((آمار = علم داده؟)) در دانشگاه میشیگان انجام داد و از آمار به عنوان علمی که به علم داده و از آماردان به عنوان افرادی که به دانشمندان داده تغییر نام داده‌اند یاد کرد.

در آپریل سال ۲۰۰۲، کمیته داده برای دانش و فناوری (CODATA)، انتشار مجله‌ای با عنوان علم داده (Data Science Journal) را آغاز کرد. این اثر، بر مسائلی مانند توصیف سیستم‌های داده، نشر آن‌ها در اینترنت، کاربردها و مسائل قانونی مربوط به این حوزه متمرکز بود.





در حدود سال ۲۰۰۷، «جیم گری» (Jim Gray)، برنده جایزه تورینگ، علوم داده محور را به عنوان چهارمین پارادایم علم معرفی کرد که از تحلیل محاسباتی داده‌های بزرگ به عنوان روشی علمی جهت ساخت دنیایی که در آن ادبیات علم و همه داده‌های علمی آنلاین هستن استفاده می‌کند.



در سال ۲۰۱۲، دانراجی پاتیل (Dhanurjay Patil “DJ”)، در مقاله «دانشمند داده: جذاب‌ترین شغل قرن ۲۱» که در مجله بررسی کسب‌وکار هاروارد منتشر شد، از دانشمندان داده به عنوان نژادی جدید یاد کرده که کمبود آن‌ها منجر به محدودیت‌های جدی در برخی از بخش‌های صنعت و دانشگاه می‌شود.

✓در سال ۲۰۱۳، اتحادیه اروپایی تحلیل داده (EuADS) در
لوکزامبورگ تاسیس شد.

✓سال ۲۰۱۵، «ژورنال بین‌المللی علم و تحلیل داده»
(International Journal of Data Science and
Analytics) توسط اسپرینگر به منظور انتشار کارهای انجام
پذیرفته در حوزه علم داده و تحلیل داده‌های کلان (مه داده)
بنا شد.



دانشمند داده به

چه کسی گفته

می شود 

جف هامربرجر (Jeff Hammerbacher) از لیدرهای

تیم تحلیل داده فیسبوک می نویسد:

یک دانشمند داده در تیم باید بتواند خط لوله پردازش چند مرحله ای در پایتون را بنویسد، یک آزمون فرضیه را طراحی کند، تجزیه و تحلیل رگرسیون را بر روی نمونه های داده با R انجام دهد، یک الگوریتم را برای برخی از محصولات یا سرویس ها در هادوپ طراحی و پیاده سازی کند یا اینکه نتایج تحلیل های ما را به دیگر اعضای سازمان انتقال دهد.



مهارت های اساسی علم داده

1. تجزیه و تحلیل آماری: الگوهای موجود در داده ها را مشخص کنید. این شامل داشتن حس شدید تشخیص الگوی و تشخیص ناهنجاری است

2. یادگیری ماشین: الگوریتم ها و مدل های آماری را پیاده سازی کنید تا رایانه بتواند به طور خودکار از داده ها یاد بگیرد.

3. علوم کامپیوتر: از اصول هوش مصنوعی ، سیستم های پایگاه داده ، تعامل انسان / کامپیوتر ، تجزیه و تحلیل عددی و مهندسی نرم افزار استفاده کنید.

✓ **برنامه نویسی:** برنامه های رایانه ای را بنویسد و مجموعه داده های بزرگ را برای کشف پاسخ به مشکلات پیچیده آنالیز کند. دانشمندان داده باید بتوانند به راحتی به زبانهای مختلفی مانند **Python، R، Java** و **SQL** کدنویسی کنند.

✓ **داستان پردازی داده ها:** با استفاده از داده ها، اغلب برای مخاطبان غیر فنی، بینش های عملی را برقرار کند.



موقعیت‌های

شغلی در حوزه

علم داده

۱. دانشمند داده (DATA SCIENTIST)

✓ متوسط درآمد سالانه: \$139,840

✓ الزامات شغلی: یافتن، تمیز کردن و سازماندهی داده‌ها برای شرکت‌ها. دانشمندان داده باید بتوانند مقادیر زیادی از اطلاعات پیچیده و خام را تجزیه و تحلیل کنند تا الگویی را پیدا کنند که به نفع یک سازمان باشد و به تصمیم‌گیری استراتژیک در زمینه تجارت کمک کند. در مقایسه با تحلیلگران داده، دانشمندان داده بسیار فنی‌تر هستند.

DATA SCIENTIST

"AS RARE AS UNICORNS"

Languages

R, SAS, Python, Matlab, SQL,
Hive, Pig, Spark

Skills & Talents

- ✓ Distributed computing
- ✓ Predictive modeling
- ✓ Story-telling and visualizing
- ✓ Math, Stats, Machine Learning



Role

Cleans, massages and organizes
(big) data

Mindset

Curious data wizard

HIRED BY



۲. مهندس یادگیری ماشین (MACHINE LEARNING ENGINEER)

✓ متوسط درآمد سالانه: ۱۱۴،۸۲۶ دلار

✓ الزامات شغلی: برنامه‌نویسانی هستند ، که تمرکز آنها بر برنامه‌نویسی هوش مصنوعی برای ماشین است. آنها برنامه‌هایی را ایجاد می‌کنند که ماشین‌ها را قادر می‌سازد به طور هوشمند برای انجام برخی از کارها اقداماتی انجام دهند. نمونه ای از سیستمی که یک مهندس یادگیری ماشین روی آن کار می‌کند رانندگی یک اتومبیل بدون سرنشین است. آنها به طور معمول به آمار و مهارت‌های برنامه‌نویسی قوی و همچنین دانش مهندسی نرم افزار نیاز دارند.

۴. معمار داده (DATA ARCHITECT)

✓ متوسط درآمد سالانه: ۱۰۸،۲۷۸ دلار

✓ الزامات شغلی: معماران داده چگونگی ذخیره، بازیابی، یکپارچه‌سازی و مدیریت داده‌ها توسط اشخاص و سیستم‌های **IT** مختلف را تعریف می‌کنند. همچنین وظیفه مدیریت برنامه‌های کاربردی را که این داده‌ها را به نوعی استفاده یا پردازش می‌کنند، دارند.

DATA ARCHITECT

THE CONTEMPORARY DATA MODELLER

Languages

SQL, XML, Hive, Pig, Spark

Skills & Talents

- ✓ Data warehousing solutions
- ✓ In-depth knowledge of database architecture
- ✓ Extraction Transformation and Load (ETL), spreadsheet and BI tools
- ✓ Data modeling
- ✓ Systems development



Role:

Creates blueprints for data management systems to integrate, centralize, protect and maintain data sources

Mindset:

Inquiring ninja with a love for data architecture design patterns

HIRED BY

VISA *Coca-Cola* logitech

۵. مهندس داده (DATA ENGINEER)

✓ متوسط درآمد سالانه: ۱۰۲,۸۶۴ دلار

✓ الزامات شغلی: پردازش دسته‌ای (**batch**) یا پردازش زمان-واقعی (**real-time**) روی داده‌های جمع‌آوری شده و ذخیره شده. مهندسان داده همچنین مسئول ساخت و نگهداری خطوط لوله داده (**data pipelines**) هستند که یک اکوسیستم داده قوی و به هم پیوسته در سازمان ایجاد می‌کنند و اطلاعات را برای دانشمندان داده در

30 دسترس می‌کنند.

DATA ENGINEER

SOFTWARE ENGINEERS BY TRADE

Role

Develops, constructs, tests and maintains architectures (such as databases and large-scale processing systems)

Mindset

All-purpose everyman



HIRED BY



Languages

SQL, Hive, Pig, R, Matlab, SAS, SPSS, Python, Java, Ruby, C++, Perl

Skills & Talents

- ✓ *Database systems (SQL & NO SQL based)*
- ✓ *Data modeling & ETL tools*
- ✓ *Data APIs*
- ✓ *Data warehousing solutions*

۶. توسعه دهنده هوش کسب و کار (BI) BUSINESS INTELLIGENCE DEVELOPER

✓ متوسط درآمد سالانه: ۸۱،۵۱۴ دلار

✓ الزامات شغلی: توسعه دهندگان BI استراتژی‌هایی را برای کمک به کاربران مشاغل در یافتن سریع اطلاعات مورد نیاز برای تصمیم‌گیری در مورد شغل بهتر، طراحی و توسعه می‌دهند. با درک بسیار بالا از داده‌ها، ابزارهای BI استفاده می‌کنند یا برنامه‌های تحلیلی BI را توسعه می‌دهند تا درک کاربران نهایی از سیستم خود را تسهیل کنند.

BUSINESS ANALYST

CHANGE AGENT

Languages
SQL

Skills & Talents

- ✓ Basic tools (e.g. MS Office)
- ✓ Data visualization tools (e.g. Tableau)
- ✓ Conscious listening and storytelling
- ✓ Business Intelligence understanding
- ✓ Data modeling



Role
Improves business processes as intermediary between business and IT

Mindset
Resilient project juggler

HIRED BY
UBER  ORACLE

۷. آماردان (STATISTICIAN)

✓ متوسط درآمد سالانه: ۷۶،۸۸۴ دلار

✓ الزامات شغلی: آماردان به منظور شناسایی روندها و روابطی که می تواند برای اطلاع رسانی در تصمیم گیری سازمانی مورد استفاده قرار گیرد، جمع آوری، تجزیه و تحلیل و تفسیر داده ها را انجام می دهند. علاوه بر این، مسئولیت های روزمره آماردان اغلب شامل فرآیندهای جمع آوری داده های طراحی، ارتباط یافتن به ذینفعان و مشاوره استراتژی سازمانی است.

STATISTICIAN

'HISTORIC LEADERS OF DATA'

Languages

R, SAS, SPSS, Matlab, Stata, Python,
Perl, Hive, Pig, Spark, SQL

Skills & Talents

- ✓ Statistical theories & methodology
- ✓ Data mining & machine learning
- ✓ Distributed Computing (Hadoop)
- ✓ Database systems (SQL and NO SQL based)
- ✓ Cloud tools



Role

*Collects, analyzes and interprets-
qualitative as well as quantitative
data with statistical theories and
methods*

Mindset

*Logical and enthusiastic stats
genius*

HIRED BY

Linked in

Johnson-Johnson



PEPSICO

۸. تحلیلگر داده (DATA ANALYST)

✓متوسط درآمد سالانه: ۶۲۴۵۳ دلار

✓الزامات شغلی: تحلیلگر داده بایستی مجموعه داده های بزرگ را متناسب با تجزیه و تحلیل مورد نظر برای شرکتها تبدیل و دستکاری کند. برای بسیاری از شرکتها ، این نقش همچنین می تواند شامل ردیابی تجزیه و تحلیل باشد. تحلیلگران داده همچنین با تهیه گزارشهایی برای رهبران سازمانی که به طور مؤثر روندها و بینشهای حاصل از تحلیل آنها را بطور مؤثر ارتباط می دهند، به فرآیند تصمیم گیری کمک می کنند.

DATA ANALYST

DATA DETECTIVE

Role

Collects, processes and performs statistical data analyses

Mindset

Intuitive data junkie with high "figure-it-out" quotient



Languages

R, Python, HTML, Javascript, C/C++, SQL

Skills & Talents

- ✓ *Spreadsheet tools (e.g. Excel)*
- ✓ *Database systems (SQL and NO SQL based)*
- ✓ *Communication & visualization*
- ✓ *Math, Stats, Machine Learning*

HIRED BY



برنامه جامع برای تبدیل شدن به دانشمند داده

- The most comprehensive Data Science learning plan for 2017

علم داده در استارت‌آپ‌ها



اولین قدم برای ساختن هر محصول داده‌ای، جمع آوری داده است. برای جمع آوری داده‌ها، ابتدا باید تصویری از کاربران پایه و تعداد لاگ رویدادهایی که به برنامه شما دسترسی دارند، باشد.

مثال: فرض کنید قصد دارید یک محصول نرم‌افزاری را بفروشید. شما نیاز دارید اطلاعاتی پایه درباره کاربران محصول خود در مراحل مختلف نصب، استفاده و خرید داشته باشید.

برای انجام این کار می‌توان برنامه را با ردیاب‌ها (**Trackers**) جاسازی کرد. قرار دادن رویدادهای ردیابی (**Tracking Events**) در سمت-کاربر (Client-side) ضروری است زیرا داده‌هایی را به سرور ارسال می‌کند که برای توسعه محصولات کاربرد دارد.

انواع داده‌های یک استارت‌آپ

Raw Data

- Do not have any schema
- Not present in any designated format
- Example: tracking events

Processed Data

- Implementation of schemas over the raw data
- Encoded in specified formats

Cooked Data

- Summary of the processed data

نکاتی درباره راه اندازی یک استارت‌آپ در حوزه علوم داده

✓ کاربران را در هر فاز پروژه همراه خود کنید.

✓ برای جلوگیری از کار مضاعف، مهندسان داده را با دانشمندان داده - حتی در مرحله نمونه سازی - درگیر کنید

✓ بیشتر جلسات به اشتراک گذاری دانش داخلی را انجام دهید و واژگان و فرایندها را تنظیم کنید.

✓ پیچیدگی سیستم را دست کم نگیرید و در مورد تمام جزئیاتی که برای تولید محصول لازم است، با کاربر صحبت کنید.

ساختن خط لوله داده‌ها (DATA PIPELINES)

پس از جمع آوری داده‌ها، شما نیاز دارید که نتایج را در کمترین زمان برای کاربران تجزیه و تحلیل، پردازش و خروجی کنید. پردازش داده‌ها مهمترین بخشی است که شرکت‌های نوپا باید از آن مراقبت کنند. بنابراین، یک خط لوله داده وظیفه پردازش داده‌های جمع آوری شده و کمک به دانشمندان داده را برای تجزیه و تحلیل داده‌ها را دارد. خط لوله داده به یک پایگاه داده وصل می‌شود که می‌تواند یک پلتفرم Hadoop یا یک SQL باشد.

تجزیه و تحلیل کیفیت محصول

تبدیل داده های خام به داده های پخته شده به نحوی که کیفیت محصول تولیدی را به شکلی مختصر و مفید گزارش دهد، یک گام مهم و اساسی برای دانشمندان داده در یک استارتاپ است. این کار می تواند توسط ابزارهایی چون تعریف KPI ها، بکارگیری زبان های برنامه نویسی قوی در گزارشگیری از داده ها مثل R و بکارگیری ETL های مختلف برای استخراج و تبدیل داده ها صورت پذیرد.

استارت‌آپ‌های موفق در حوزه علم داده



مای اسمارت ژن



مای اسمارت ژن استارت‌آپی در حوزه تحلیل داده‌های ژنتیکی است که مأموریت آن کشف استعدادها و خاصیت هر فرد است. استارت‌آپ مای اسمارت ژن توسط منصور حسنی (کارشناسی ارشد هوش مصنوعی) و دکتر دامون نشتاعلی حسنی (دکتری هوش مصنوعی) در بهمن ۹۷ بنیان‌گذاری شده است.

استعدادیابی ژنتیکی با نمونه بزاق دهان

بهترین تصمیم‌ها وقتی گرفته می‌شوند که شناخت بهتری از خودمان داشته باشیم
شناختنامه ژنتیکی راهنمای مخصوص شما در این مسیر است

📞 راهنمایی آنلاین در واتساپ

دیتابین

دیتابین یک سامانه تحلیل رفتار کاربر برای اپلیکیشن‌های موبایلی است. توسعه‌دهندگان می‌توانند SDK دیتابین را در برنامه خود پیاده‌سازی کنند تا داده‌های هر کاربری که برنامه آنها را نصب می‌کند، برای دیتابین ارسال شود. ندا مثنوی (ارشد کارآفرینی) بنیانگذار این استارت‌آپ است که به همراه مهدی یزدی (کارشناس فناوری اطلاعات) و حامد پالیک (کارشناس کامپیوتر) آن را راه‌انداخته‌اند.



پایشگر هشتگ (سامانه پایش شبکه های اجتماعی)

این محصول به صورت آنلاین و لحظه ای، آخرین اخبار مرتبط با شما را از طریق پایش فضای مجازی در قالب داشبوردی اختصاصی به مدیران و کارشناسان روابط عمومی، بازاریابی، فروش و واحدهای پشتیبانی و ارتباط با مشتریان ارائه می نماید.

پایش، تحلیل و جستجو در شبکه های اجتماعی

بروزترین و جامع ترین سامانه پایش شبکه های اجتماعی را امتحان کنید!



NAUTO استارتاپ دَش کمرای هوش مصنوعی برای وسایل نقلیه

دَش کمرای نوعی دوربین مخصوص نصب در بخش مقابل راننده یا داخل و شیشه جلوی اتومبیل است که وقتی اتومبیل در حال حرکت است، این دوربین به طور پیوسته از راه جلوی اتومبیل تصویربرداری می کند. این محصول تحت شبکه، مبتنی بر فناوری ابری است که خطرات تهدیدکنندهی رانندگان را تشخیص می دهد. این استارتاپ برای دریافت بازخورد در آخر سفر و آنالیز دلایل تصادفات امکانی فراهم کرده است.



REFLEKTIVE استارتاپ

بازخورد مداوم عملکرد کارمندان

رفلکتیو، نرم‌افزار دریافت و ثبت گزارش و بازخورد عملکرد به کارکنان است. این نرم‌افزار امکانی برای مدیران فراهم می‌کند که از اهداف و شیوه‌ی عملکرد کارکنان چک‌لیستی آنی داشته باشند و بتوانند دربارهی هر کدام از آنها، گزارش مستقیمی ارائه دهند. برای کارکنان هم، این نرم‌افزار نوعی کارنامه‌ی در حال اجراست که آنها را هر لحظه از نحوه‌ی عملکردشان و میزان رضایت مدیران از عملکردشان مطلع می‌کند.

AYLIEN استارتاپ

پلتفرم هوشمند اخبار

آیلین یک پلتفرم علوم داده در ایرلند است که راه حل‌های تحلیل محتوا با قدرت AI را ارائه می‌دهد. مدیرعامل این استارتاپ یک ایرانی به نام پارسا غفاری است. ایلین بسته‌ای از هشت ابزار برای پردازش زبان طبیعی (NLP)، یادگیری ماشین و بازیابی اطلاعات است تا به راحتی معانی ارزشمندی را از اسناد استخراج کند.



دیدن استارت‌آپ‌هایی بیشتر با موضوع علم داده:

1. [European data science start-ups to watch](#)
2. [42 Big Data Startups](#)
3. [Deaswatch](#)

منابع:

1. Tansley, S., & Tolle, K. (2009). *The fourth paradigm: data-intensive scientific discovery* (Vol. 1). T. Hey (Ed.). Redmond, WA: Microsoft research.

2. علم داده چیست؟ مجله علمی فرادرس

پایان