

DATA to AI Lab

LIDS, MIT

Kalyan Veeramachaneni

Over the past 7-8 years

Data Scientist: The Sexiest Job of the 21st Century

by **Thomas H. Davenport** and **D.J. Patil**

FROM THE OCTOBER 2012 ISSUE

Help Wanted: Black Belts in Data

Starting salaries for data scientists have gone north of \$200,000

by **Rodrigo Orihuela** and **Dina Bass**

June 4, 2015, 1:07 PM EDT

Updated on June 4, 2015, 2:00 PM EDT

From **BloombergBusinessweek** | [Subscribe](#) | [Reprints](#)

Artificial intelligence (AI)

2016: the year AI came of age

Google and Amazon brought AI into the home and DeepMind built a computer that could outsmart humans at Go. Will 2017 hold similar advancements?

Connection between AI and Data Science?

- All applications of AI that you see have been developed by
 - Collecting data
 - Learning models from them
 - Using those models to act
 - Take for example – alphago
 - “The system's neural networks were initially bootstrapped from human gameplay expertise. AlphaGo was initially trained to mimic human play by attempting to match the moves of expert players from recorded historical games, using a database of around 30 million moves”
- Predictive analytics

And The Winner Is: Big Data Oscar Picks

This year's Best Picture will be 12 Years A Slave, according to Academy Award prognosticators at Farsite. Will big data get it right?

Academic Research

ICML @ Sydney

Thirty-fourth International
Conference on Machine
Learning

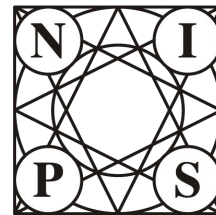
Since 1980



KDD2017

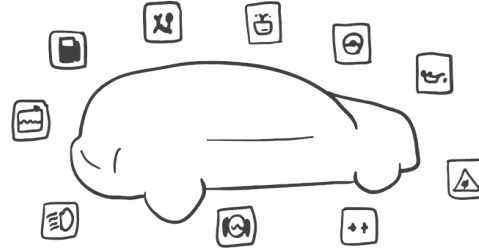
Since 1995

*Conference on Neural Information
Processing Systems*



Since 1986

Numerous Data Science Problems



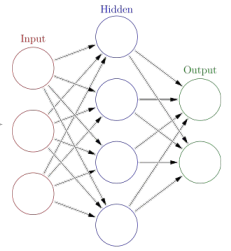
Tremendous increase in rate at which we are encountering data science problems.
The challenge is not to solve just one problem, but to overcome the bottlenecks
that prevent us from solving many!



Build AI products faster



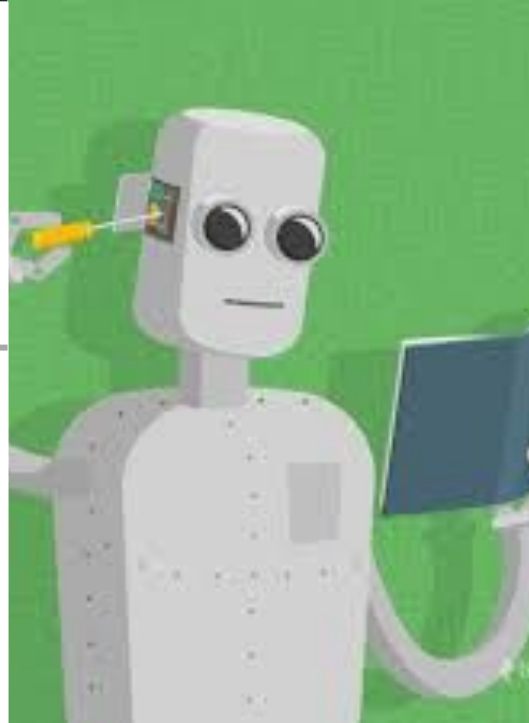
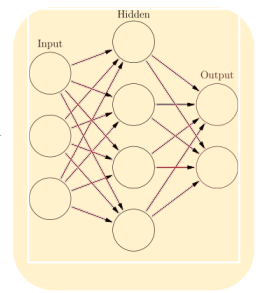
Age	Score	...	Default
32	678		y
21	786		n
67	776		n



Training

Usage

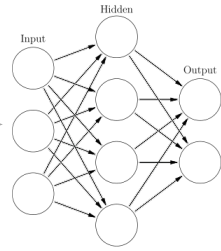
Age	Score	Default
25	524	...	?



Building AI products



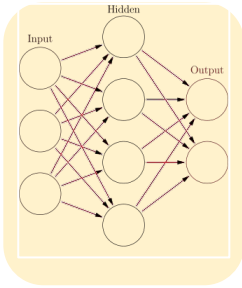
Age	Score	...	Default
32	678		y
21	786		n
67	776		n



Training

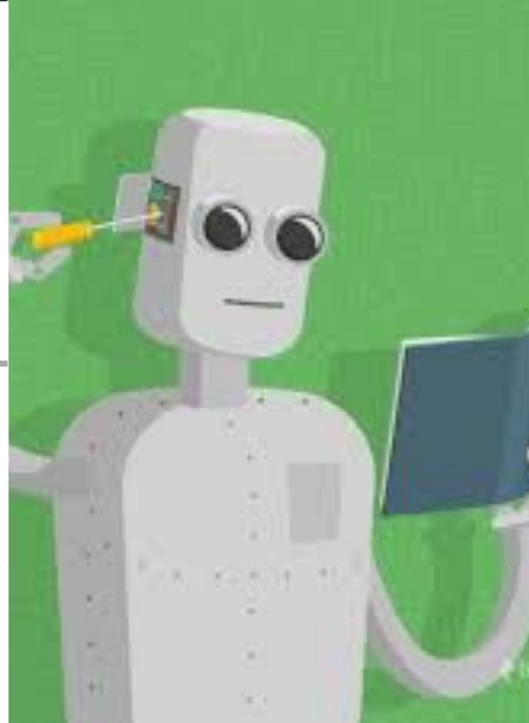
Usage

Age	Score	Default
25	524	...	?



No

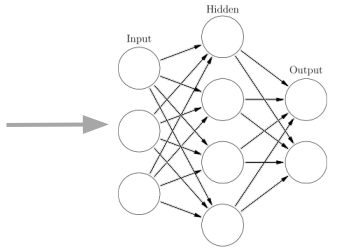
New user



Building AI products



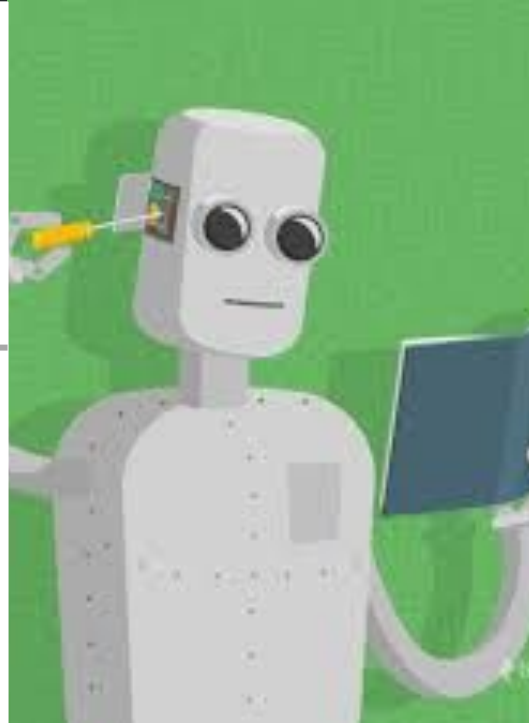
Age	Score	...	Default
32	678		y
21	786		n
67	776		n



Training

But we can train many more models for different outcomes from the same data ..

- When is a user most likely to refinance ?
- When is a user most likely to buy next car?
- And so on.....

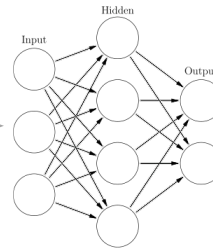


Building AI products

To unlock the potential of machine learning

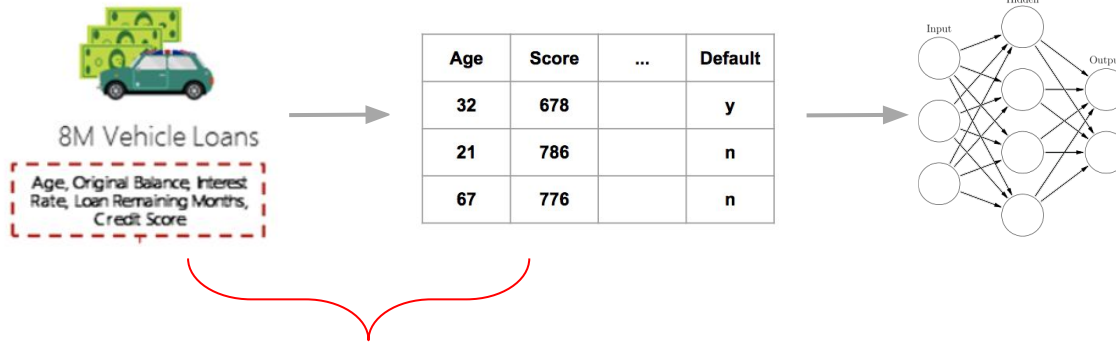


Age	Score	...	Default
32	678		y
21	786		n
67	776		n



Automatically choose and learn a model

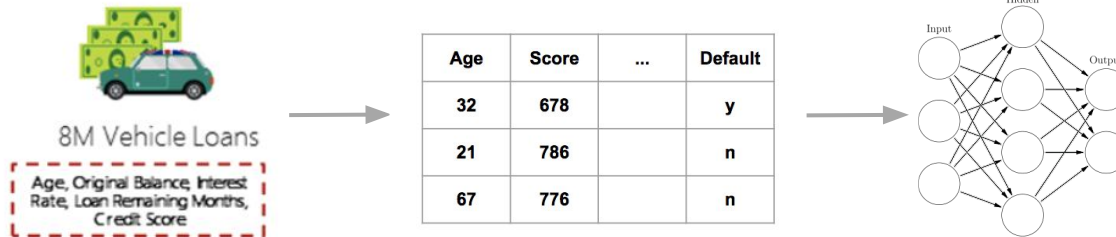
To unlock the potential of machine learning



Automatically form patterns from historical data

- Number of times customer was delayed in payments
- Rate of change of user's salary profile
- Rate of change of user's credit score

To unlock the potential of machine learning



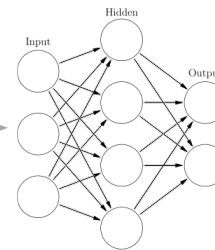
Automatically formulate questions

Change from loan default prediction to predict refinancing

To unlock the potential of machine learning

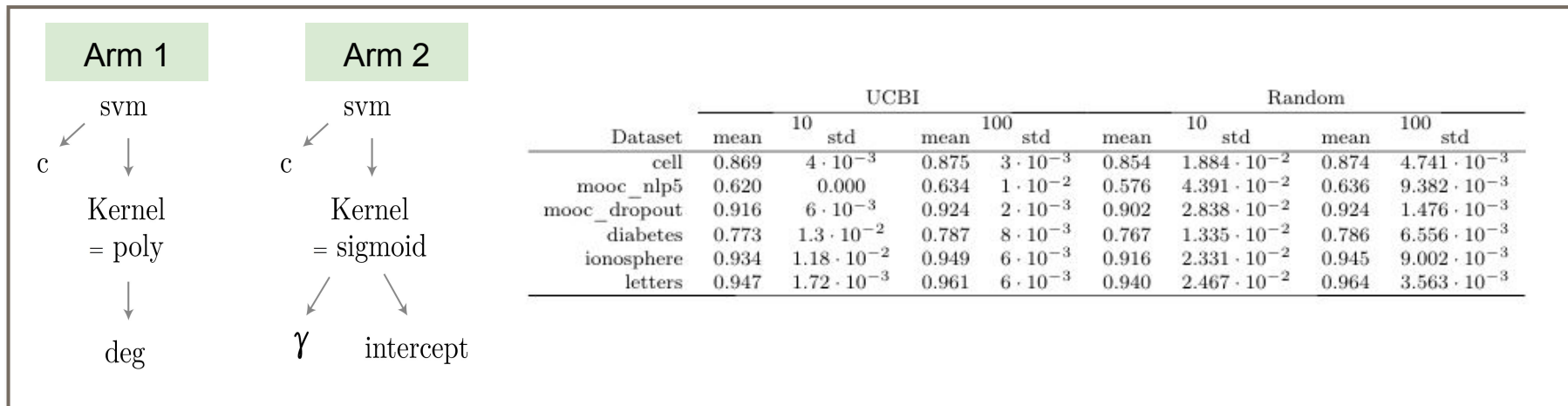


Age	Score	...	Default
32	678		y
21	786		n
67	776		n



Automatically choose and learn a model

ATM – Automatically choose a model



Pick an arm using Multi Armed bandit. Tune the hyperparameters using Gaussian Processes.
Circa 2014.

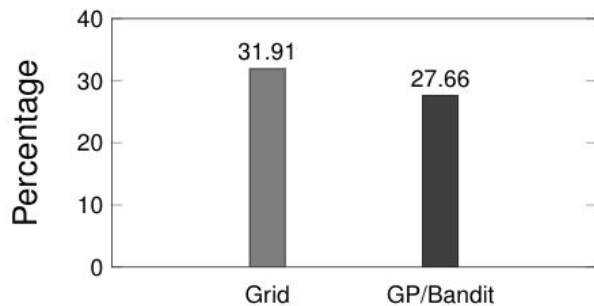
A distributed, multi-model, self-learning platform for machine learning

US Patent Application 14/598,628

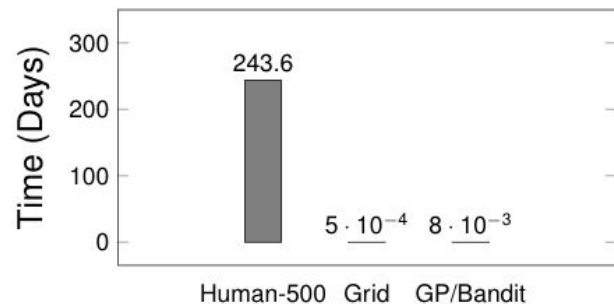
Filed January 16, 2015



ATM - Open source release and comparing to humans



Tested on 420 publicly available datasets
3 million models trained and counting
Compared against human baselines on OpenML
Ready to use!

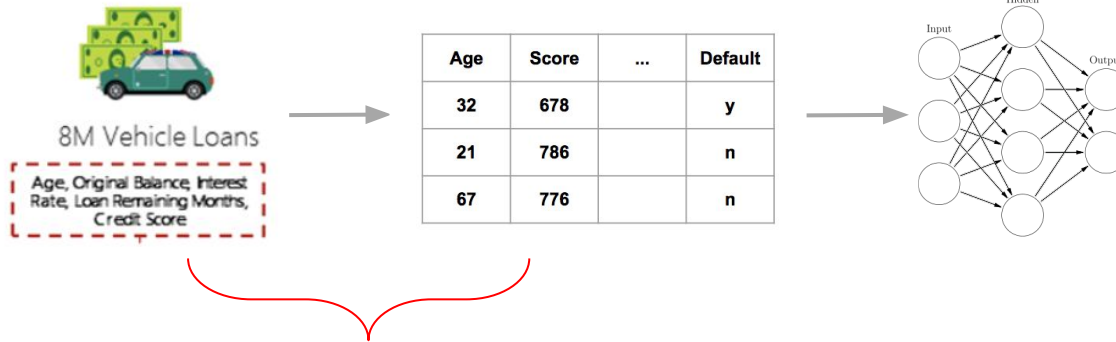




ATM

<http://bit.ly/MIT-HDI>

To unlock the potential of machine learning



Automatically form patterns from historical data

- Number of times customer was delayed in payments
- Rate of change of user's salary profile
- Rate of change of user's credit score

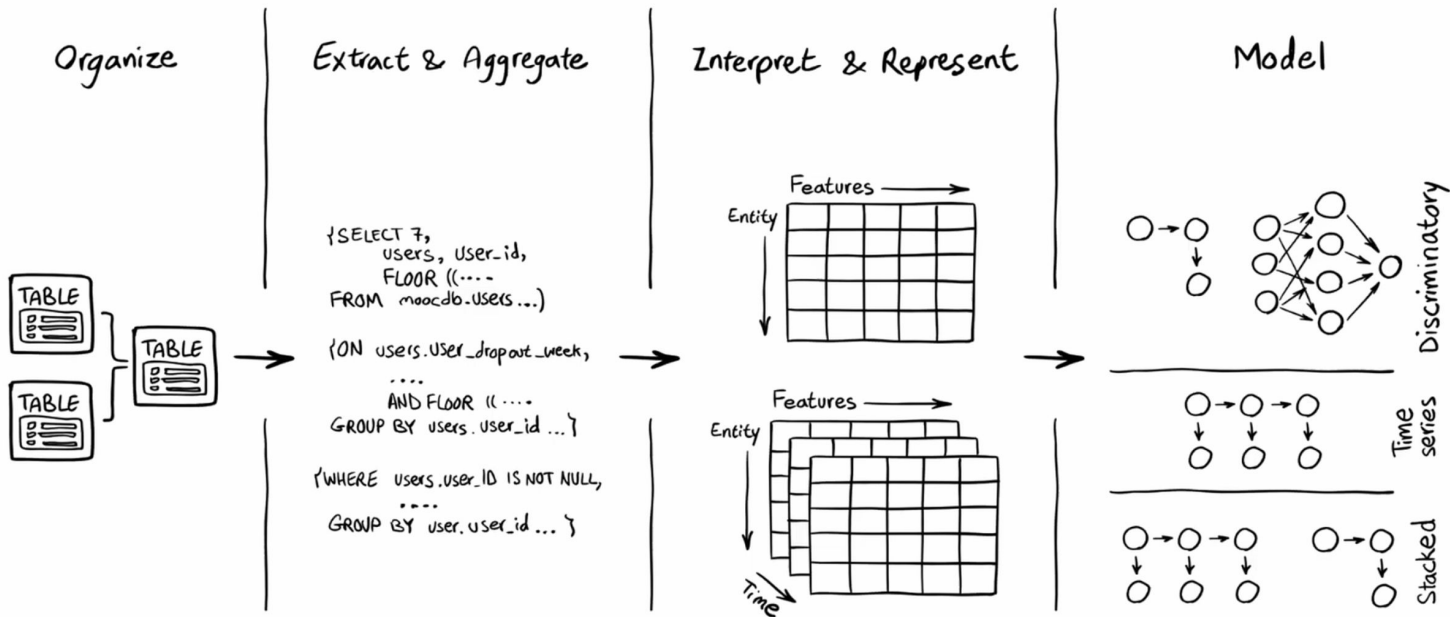
The Quintessential Matrix

FEATURIZED DATA

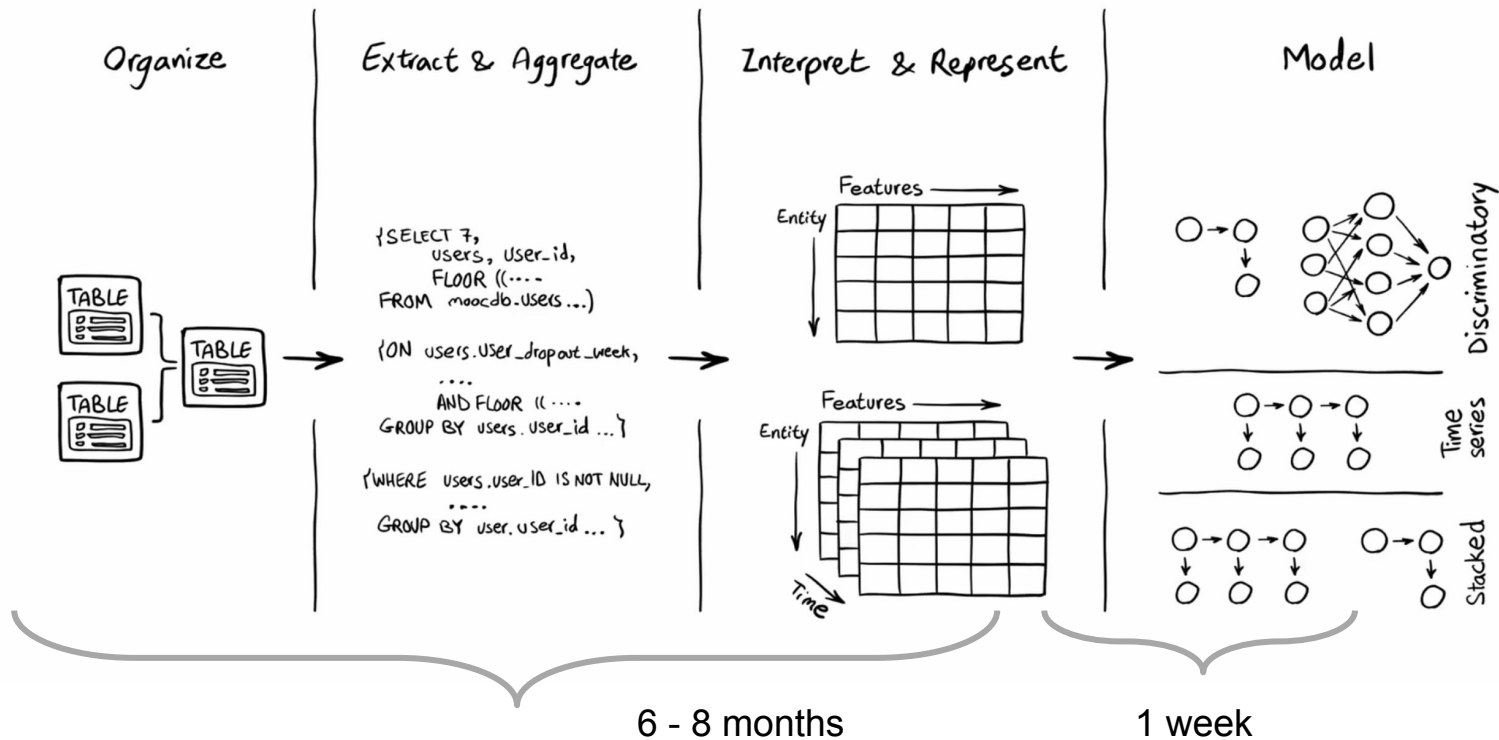
Id	X_2	X_3	X_4	X_n	Y

First column is id – project(accenture), car(Jaguar) – training example.
Each column is a feature. Last column is the label.

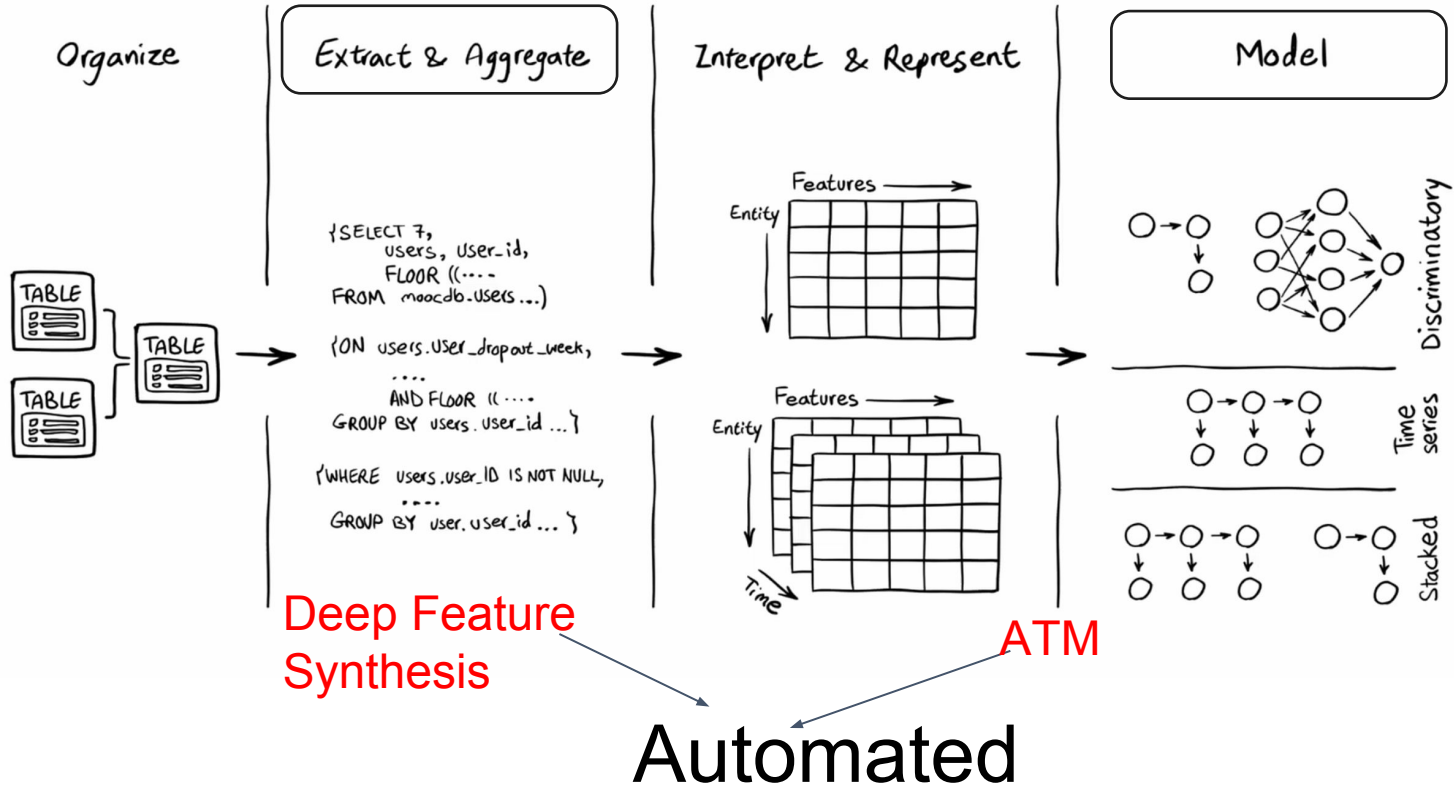
So we started the process of "Data Science"



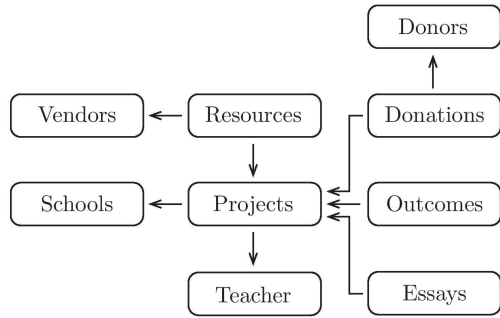
So we started the process of "Data Science"



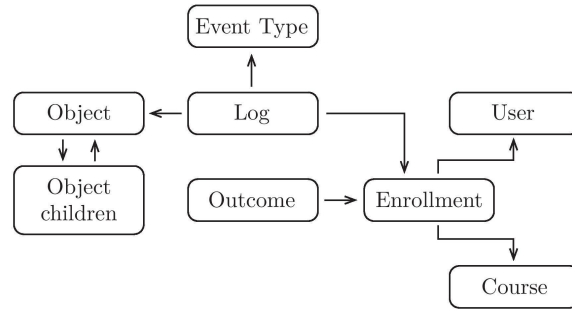
The data science process



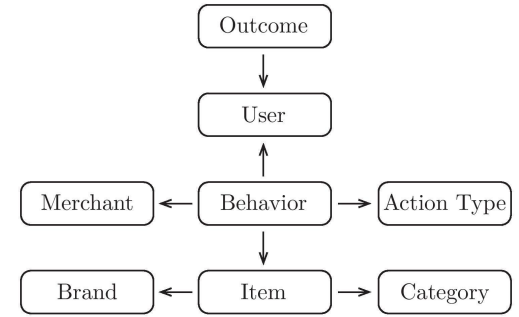
Testing our automated feature engineering on KAGGLE Competitions



Project Excitement



Dropout Prediction

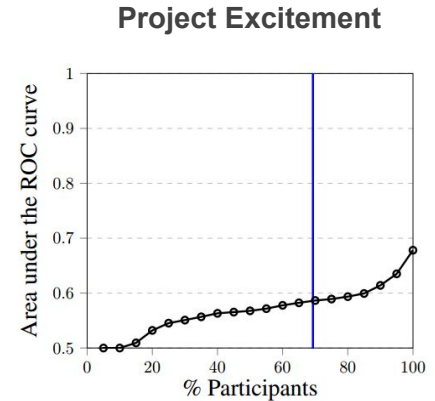
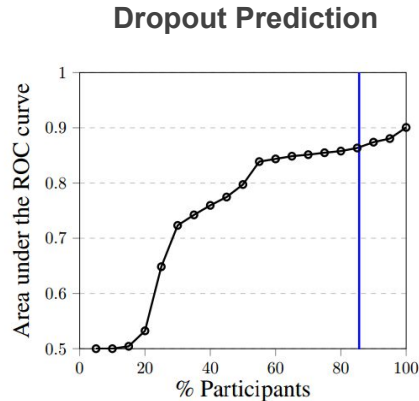
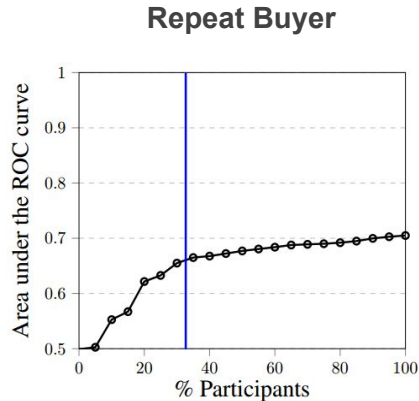


Repeat Buyer

Automation – Circa. 2015

GIVEN LEARNING SEGMENTS

We can generate features, learn models and evaluate



Lines show the standing of the Deep Feature Synthesis in the competition as of May 18th, 2015

Automation – Circa. 2015

GIVEN LEARNING SEGMENTS

We can generate features, learn models and evaluate

Tested against

1,000

data scientists

On average

92%

of top score

Over

1,200

days saved

Automation – Circa. 2015

The Washington Post
Democracy Dies in Darkness

Speaking of Science

New MIT algorithm rubs shoulders with human intuition in big data analysis

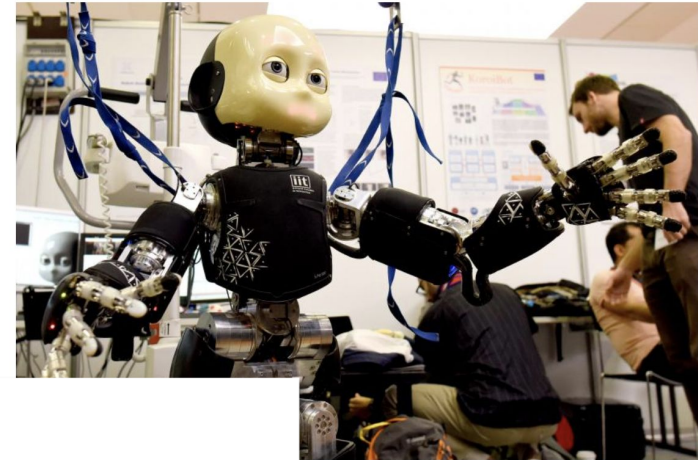
By Rachel Feltman October 19, 2015

Automating big-data analysis

System that replaces human intuition with algorithms outperforms 615 of 906 human teams.

AN ALGORITHM MAY BE BETTER THAN HUMANS AT BREAKING DOWN BIG DATA

BY SEUNG LEE ON 10/19/15 AT 6:25 PM



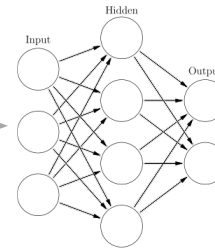


www.featuretools.com

To unlock the potential of machine learning



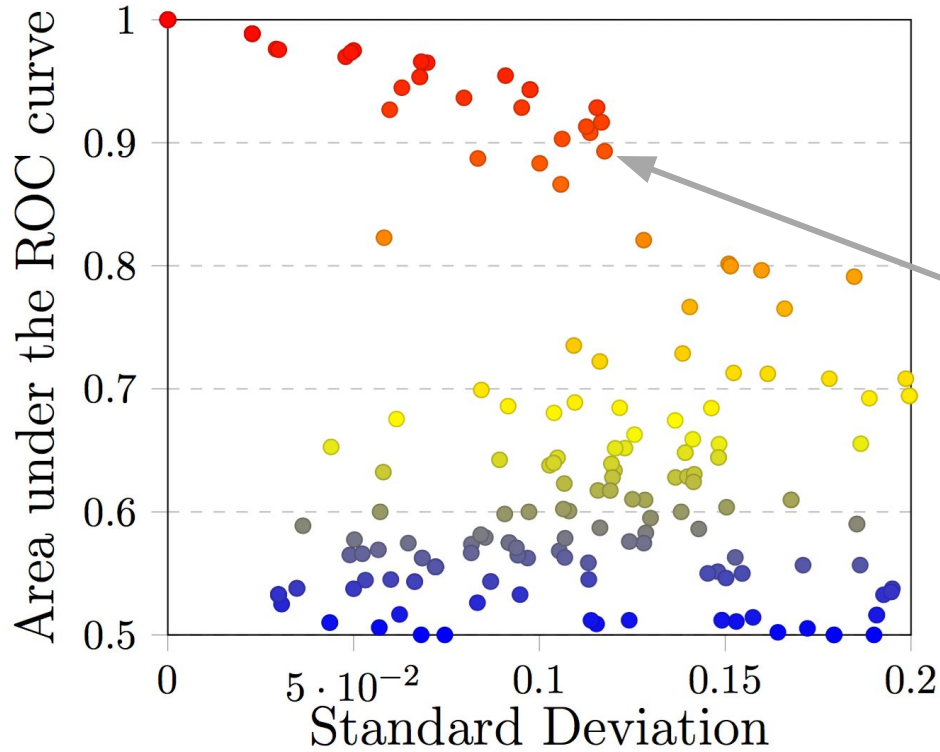
Age	Score	...	Default
32	678		y
21	786		n
67	776		n



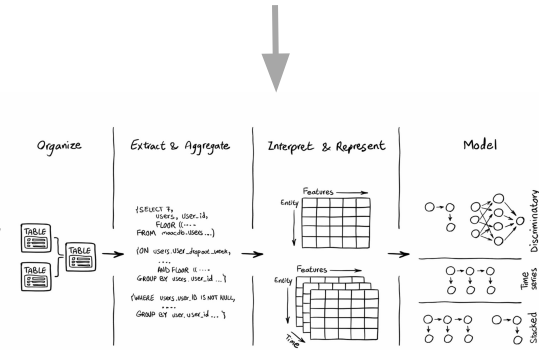
Automatically formulate questions

Change from loan default prediction to predict refinancing

TRANE



Predict whether the mean sales volume will exceed \$5000 in a 2-week window



MIT - The Human Data Interaction Project



Deep Mining
○

featurehub



SenseML

Trane

ADEL
AUTOMATIC DATA ELEMENTS LINKING

SDV
SYNTHETIC DATA VAULT

Industrial scale problems



Industrial scale problems



Predict destination
7000 fields
> 1 years data

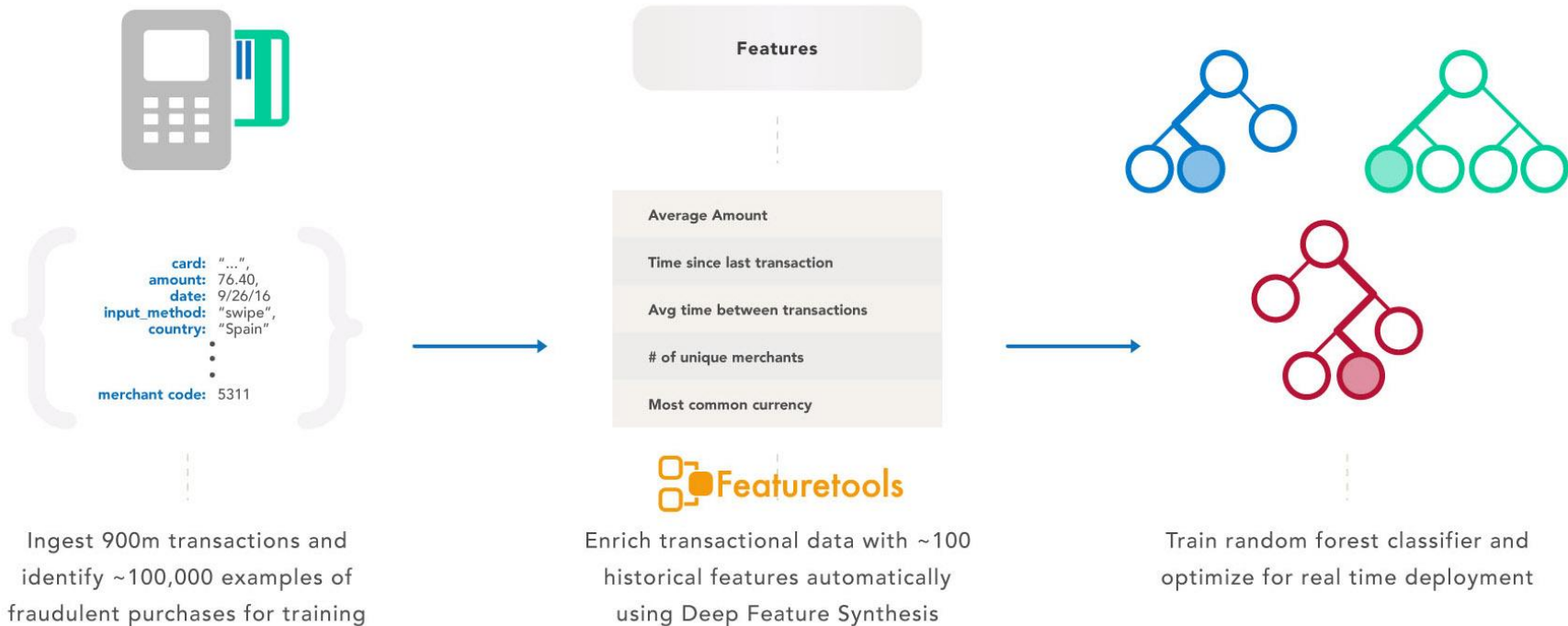


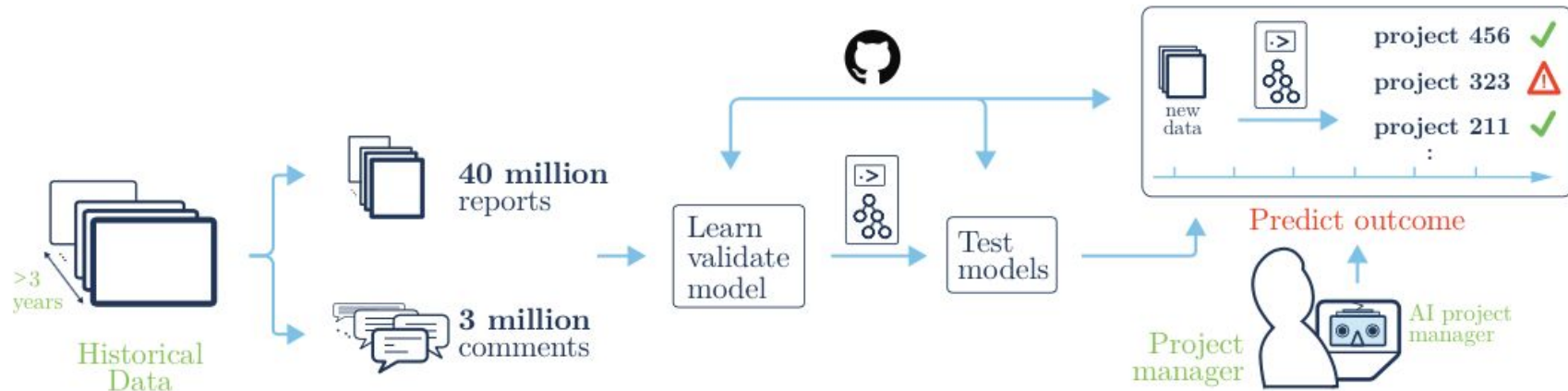
The image shows the Accenture logo, which features a yellow chevron symbol above the word "accenture" in a bold, lowercase, sans-serif font.

High performance. Delivered.

**Predicting software release
delays**

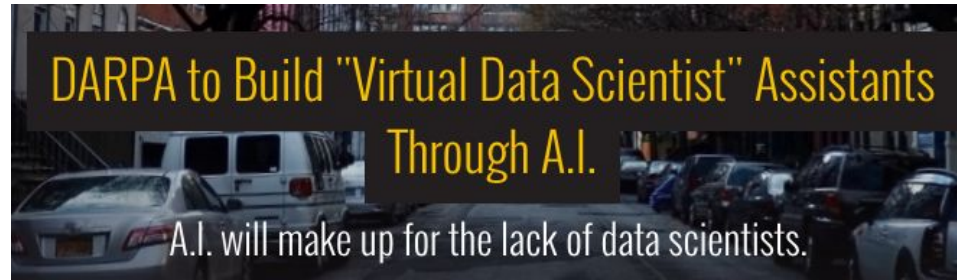
512 fields, 5 tables
>5 years data





What does the future look like?

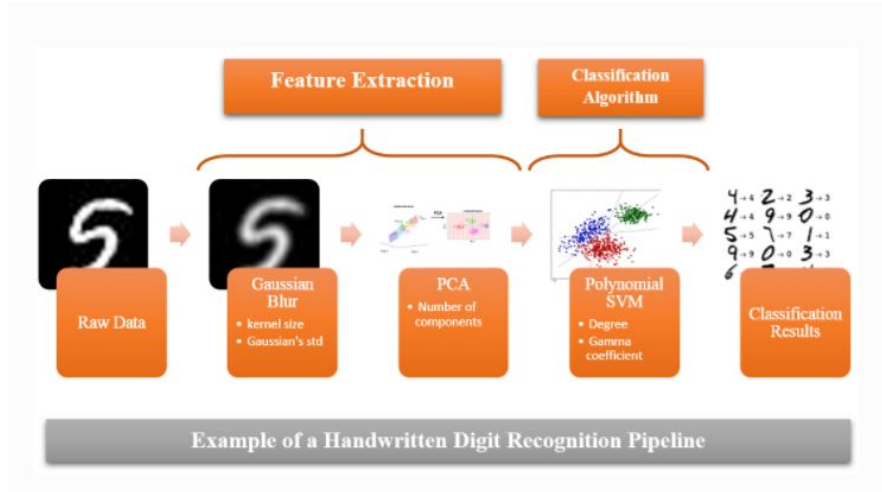
2016



Deep Mining - Much broader system

Deep Mining project aims to construct a end to end Machine Learning system automatically and for all data types.

Here is an example: The [handwritten digit recognition problem](#)



Our automated system – Feb 2018

MIT+FeatureLabs				
A	B	C	D	E
TA2 system	Coverage	% of datasets the system did best?	How many datasets the system was best?	
aika (Berkeley)	50%	11%	2 out of 18	
brown	39%	11%	2 out of 18	
columbia_u_uchicago	0%	0%	0 out of 18	
cra_eve	0%	0%	0 out of 18	
featurelabs_mit_btbt (MIT)	78%	44%	8 out of 18	
isi (USC)	33%	0%	0 out of 18	
nyu	39%	6%	1 out of 18	
qntfy	33%	11%	2 out of 18	
sri_tpot	61%	11%	2 out of 18	
texasam_tamu	94%	22%	4 out of 18	
Tests were out of 18 "seed" datasets				
Coverage --> how many datasets the system ran - trained a model, outputted a model and NIST was able to test and score the model - without an error				
Best performance --> how many datasets did the system do better than everyone else				
		#1		
		#2		

MIT - The Human Data Interaction Project

xylem
Let's Solve Water



accenture

SES[▲]
your satellite company

<http://bit.ly/MIT-HDI>



Goal: 10 new industrial scale applications!

To receive updates, apply to be one of our partner send email to:

dailabmit@gmail.com

or

kalyanv@mit.edu

