



Data Science and Data Engineering

Quick Overview

datascience**dojo**
unleash the data scientist in you



Data Science & Machine Learning

- ~3 days
- 60% theory. 40% Hands-on Exercises
- Math/Theory is minimal but not trivial
- Primary tools: R and Azure ML Studio
- Mentored Kaggle participation
- Emphasis on best practices

datascience**dojo**
unleash the data scientist in you

Data Engineering

- Teach enough data engineering skills to be effective data scientist
- 20% theory. 80% hands-on
- Handle volume, variety and velocity of data
- Internet of Things (IoT) hack day.

datascience**dojo**
unleash the data scientist in you

Hack Day

- Gather temperature and humidity data in real-time
- Use message queues, stream processors to get real time analytics



datascience**dojo**
unleash the data scientist in you

Logistics

- Eight hours of pre-bootcamp work
- Bootcamp: May 25th – May 29th (8:30 am – 5:00 pm)
- Hackday: May 30th
- Slides, sample code and other resources are consolidated in a git repository:
 - Send your github id to yuhui@datasciencedojo.com
- Office hours. Kaggle. LinkedIn group

datasciencedojo
unleash the data scientist in you

Introduction to Big Data, Predictive Analytics, and Data Science

datasciencedojo
unleash the data scientist in you

Big Data and Data Science Everywhere



data science dojo
unleash the data scientist in you

Online Shopping

Best Value

Buy **Predictive Analytics: The Power of Prediction** and **How to Measure Anything: Finding the Value of Intangibles in Business** at an **additional 5% off** Amazon.com's everyday low price.

Buy together today: \$45.43

[Add both to Cart](#)


[Show availability and shipping details](#)

Customers Who Bought This Item Also Bought


 <p>Predictive Analytics: Microsoft Excel Conrad Carlberg ★★★★☆ (10) Paperback \$24.36</p>	 <p>Big Data: A Revolution That Will Transform ... Viktor Mayer-Schönberger ★★★★☆ (32) Hardcover \$15.84</p>	 <p>Big Data: Big Analytics: Emerging Business ... Michael Minelli ★★★★☆ (6) Hardcover \$32.82</p>	 <p>How to Measure Anything: Finding the Value of ... Douglas W. Hubbard ★★★★☆ (56) Hardcover \$31.96</p>	 <p>Secrets of Analytical Leaders: Insights ... Wayne Eckerson ★★★★☆ (10) Perfect Paperback \$44.96</p>	 <p>Big Data Analytics: Disruptive ... Dr. Arvind Sathi ★★★★☆ (5) Paperback \$10.45</p>
--	--	--	---	--	---

data science dojo
unleash the data scientist in you

Social Networks









Who to follow

Twitter recommends suggested for you based on who you follow and more.

Search using a person's full name or @username



DataQualityPro.com @dataqualitypro
The most popular online data quality community resource for anyone requiring free expert tutorials, techniques, articles or technology advice.
Followed by Big Data Science and Data Science Central.




Stat Fact @StatFact
One statistics tip per day M-F from @JohnDCook. See also @ProfFact, @CompSciFact, and @SciPyTip.
Followed by Data Science Central and Big Data Science.




Anthony Goldbloom @antgoldbloom
Founder and CEO of Kaggle.


People You May Know See All




Andres Ponce



Jessica Clark
1 mutual friend



Melody Vilantino
7 mutual friends



Isabella Lopez
2 mutual friends

JOBS YOU MAY BE INTERESTED IN


TIGER Software Developer, Data Analytics
Tiger Analytics - Raleigh, NC

a Machine Learning Scientist
Amazon - Greater Seattle Area

Microsoft Data Scientist, Senior - OSD D&A
Microsoft - Bellevue, WA, US

Microsoft Data Scientist, Senior - Bing - D&A
Microsoft - Bellevue, WA, US

[Feedback](#) | [See more](#)




unleash the data scientist in you

Online Entertainment

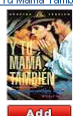
Netflix

Other Movies You Might Enjoy




Add

★★★★☆
☐ Not Interested




Add

★★★★☆
☐ Not Interested




Add

★★★★☆
☐ Not Interested




Add

★★★★☆
☐ Not Interested



Add

★★★★☆
☐ Not Interested



Add

★★★★☆
☐ Not Interested

Eiken has been added to your Queue at position 2.

This movie is available now.

[Continue Browsing](#)
[Visit your Queue](#)

data science dojo
unleash the data scientist in you

5

Brainstorming

- What are some other applications?

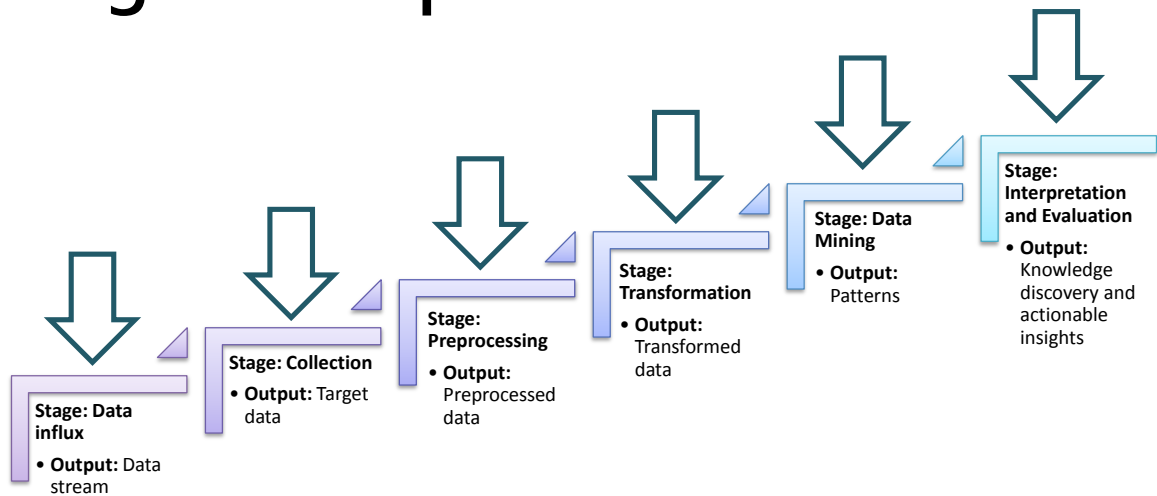
data**science**dojo
unleash the data scientist in you

Connecting the Dots

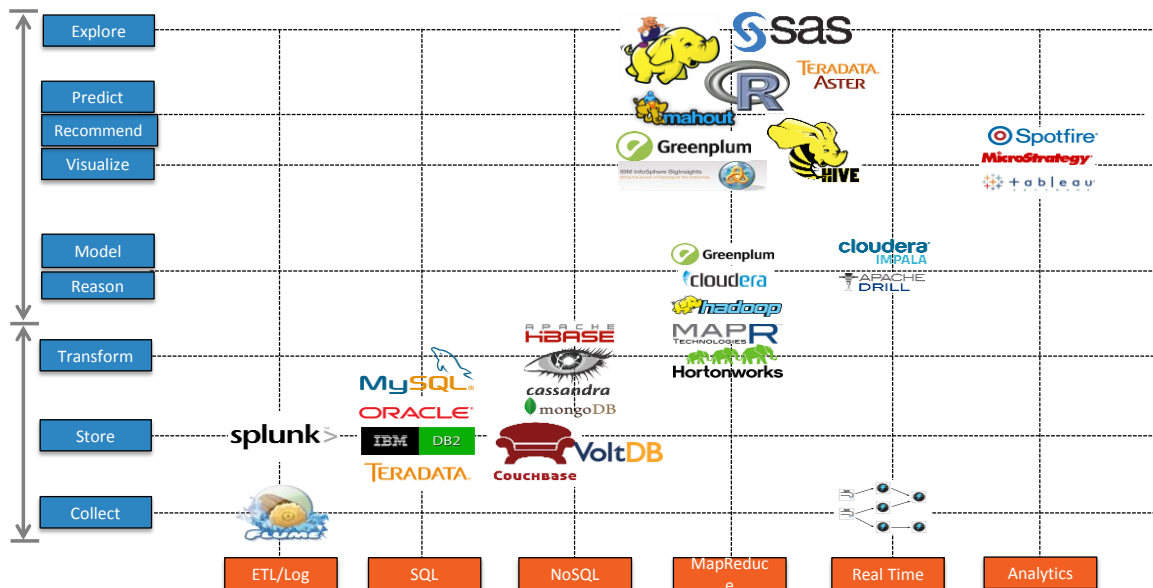
- The underlying magic behind what we saw is 'big data' and 'predictive analytics'



Big Data Pipeline



data science dojo
unleash the data scientist in you



Big Data – Technology, Platforms & Products

data science dojo
unleash the data scientist in you

Data Mining Tasks

■ Descriptive Methods:

- Find human-interpretable patterns that describe the data
- Techniques: Clustering, Association Analysis, x-point summaries

■ Predictive Methods:

- Use available data to build models that can predict the outcome of future data
- Techniques: Classification, Regression, Anomaly, and Deviation Detection

■ Prescriptive Methods:

- Predict future outcomes and suggest actions that may prevent or mitigate the impact of the predicted outcomes
- Techniques: Various optimization techniques

data science dojo
unleash the data scientist in you

Traffic Management



Descriptive [Informing Role]:

- Traffic jam has happened already.
- [Implicit: Do something about it.]

data science dojo
unleash the data scientist in you

Traffic Management



Predictive [Informing and Warning Role]:

- Traffic jam is about to happen in the next 30 minutes.
- [Implicit: Do something before it happens.]

data science dojo
unleash the data scientist in you

Traffic Management



Prescriptive [Informing, Warning, and Advisory Role]:

Take action so traffic jam does not happen

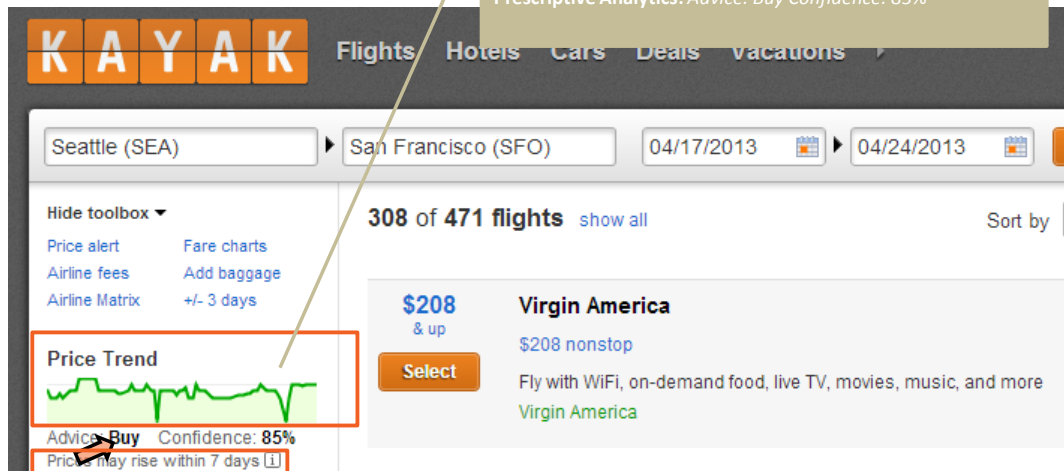
OR

Traffic jam is about to happen in the next 30 minutes and you could possibly take the following courses of action:

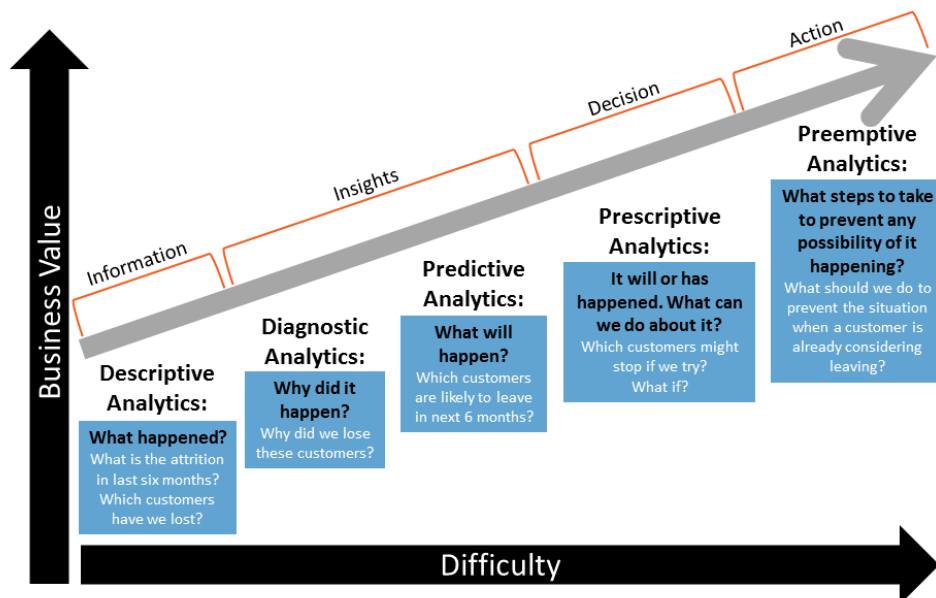
- Route traffic to service road near I-5
- Block more traffic from entering the WA-520 bridge

data science dojo
unleash the data scientist in you

Online Travel



unleash the data scientist in you



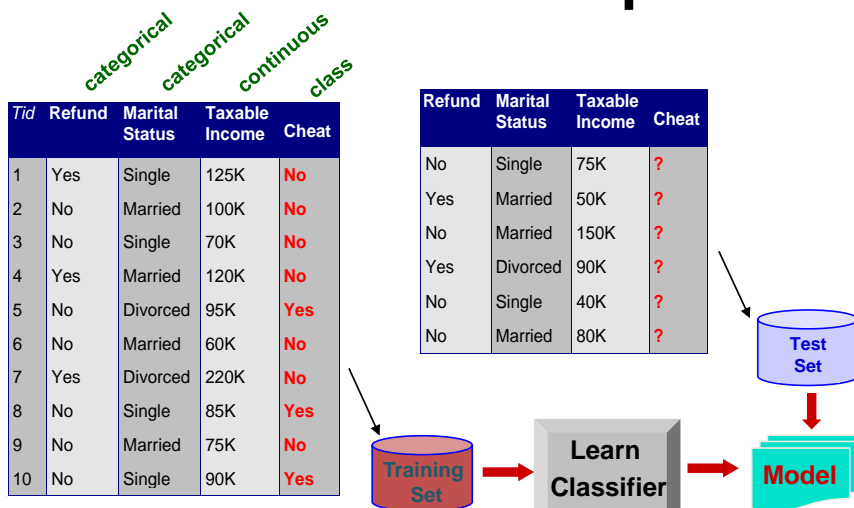
data science dojo
unleash the data scientist in you

Data Mining and Predictive Analytics

In the next few slides, we will take a look at some of the most common data mining tasks.

data**science**dojo
unleash the data scientist in you

Classification: A Simple Example



data**science**dojo
unleash the data scientist in you

Classification

- Given a collection of records (**training set**)
 - Each record contains a set of *attributes*; one of the attributes is the *class label*.
- Find a **model** for class attribute as a function of the values of other attributes.
- **Goal:** previously unseen records should be assigned a class as accurately as possible.

data**science**dojo
unleash the data scientist in you

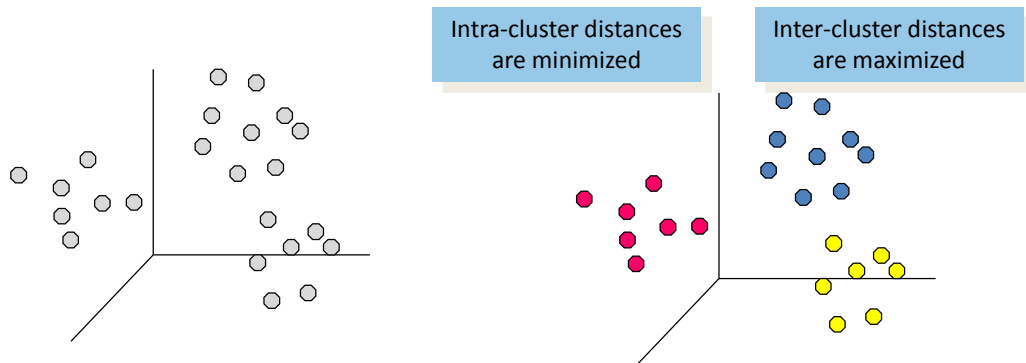
Classification: More Examples

- **Direct Marketing**
 - Goal: reduce cost of mailing by targeting a set of consumers likely to buy a new cell-phone product
- **Fraud Detection**
 - Goal: predict fraudulent cases in credit card transactions
- **Customer Attrition/Churn**
 - Goal: predict whether a customer is likely to be lost to a competitor

data**science**dojo
unleash the data scientist in you

Clustering: An Illustration

Clustering in 3-D space using Euclidean distance



data**science**dojo
unleash the data scientist in you

Clustering: Examples

- Subdivide the market into distinct subsets of customers where any subset may conceivably be selected as a segment to be reached with a particular offer



Clustering

- Given a set of data points, each having a set of attributes, and a similarity measure among them, find clusters such that:
 - Data points within a cluster have more similarities with one another
 - Data points in different clusters have less similarities with one another

data**science**dojo
unleash the data scientist in you

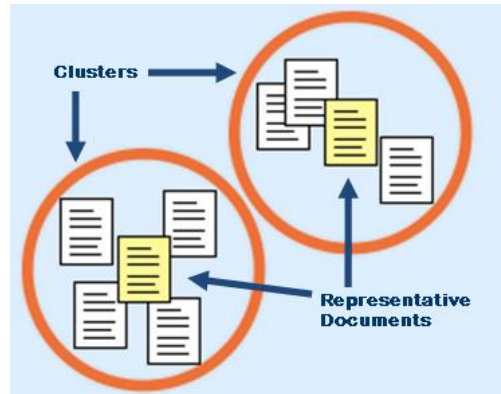
Clustering: Similarity Measures

- **Similarity Measures:**
 - Euclidean Distance if attributes are continuous
 - Other problem-specific measures
 - **Example:** If a particular word occurs in two documents or not

data**science**dojo
unleash the data scientist in you

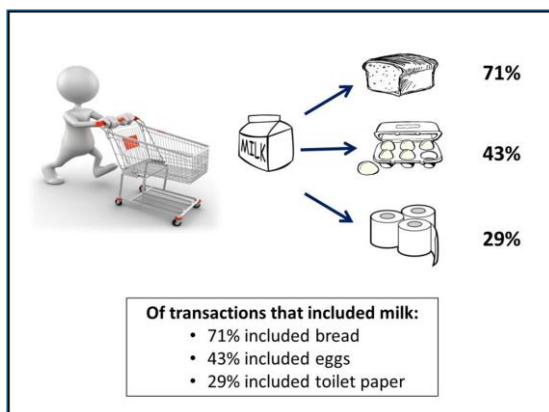
Clustering: Examples

- To find groups of documents that are similar to each other based on the important terms appearing in them



datascience**dojo**
unleash the data scientist in you

Association Analysis



Your behavior is
being predicted,
not by studying
you, but by
studying others.

datascience**dojo**
unleash the data scientist in you

Association Rule Discovery

- Given a set of records each of which contain some number of items from a given collection:
 - Produce dependency rules which will predict the occurrence of an item based on the occurrences of other items

TID	Items
1	Bread, Coke, Milk
2	Beer, Bread
3	Beer, Coke, Diaper, Milk
4	Beer, Bread, Diaper, Milk
5	Coke, Diaper, Milk

Rules Discovered:

$\{\text{Milk}\} \rightarrow \{\text{Coke}\}$
 $\{\text{Diaper, Milk}\} \rightarrow \{\text{Beer}\}$

data**science**dojo
 unleash the data scientist in you

Association Analysis: Supermarket Shelf Management

- Goal: To identify items that are bought together by a sufficient amount of customers
- Place the items close to each other on supermarket shelves



data**science**dojo
 unleash the data scientist in you

Association analysis examples

- **Marketing and sales promotion:**
 - Users who buy item A usually also buy item B
 - If users bought item A, suggest item B or even offer discount on item B
- **Inventory management:**
 - Goal: A consumer appliance repair company wants to anticipate the nature of repairs on its consumer products and keep the service vehicles equipped with the right parts to reduce the number of visits to consumer households

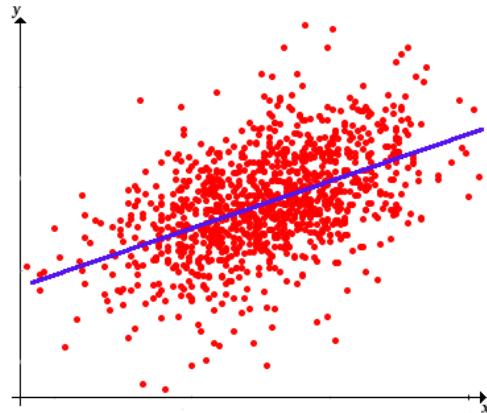
datascience**dojo**
unleash the data scientist in you

Regression Example: Predict Housing Prices



Regression

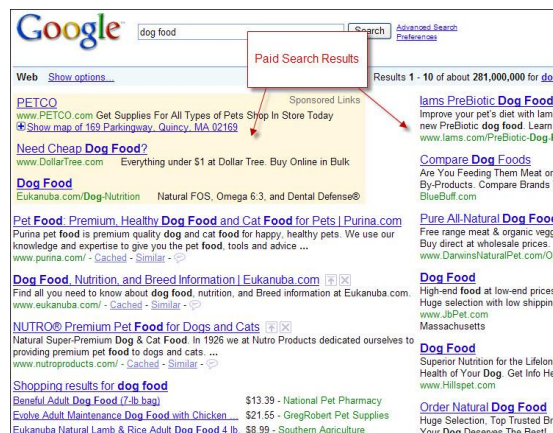
- Predict a value of a given continuous valued variable based on the values of other variables, assuming a linear or nonlinear model of dependency



data**science**dojo
unleash the data scientist in you

Regression: Ad Clicks

Predict the probability of whether or not an ad will be clicked



data**science**dojo
unleash the data scientist in you

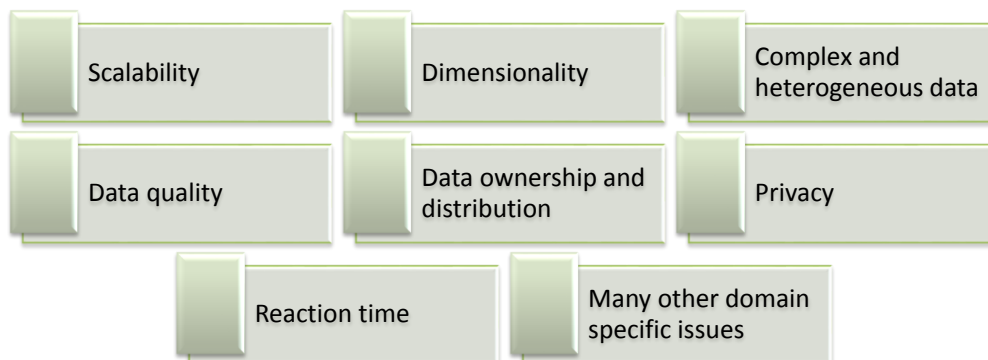
Deviation/Anomaly Detection

- Detect significant deviations from normal behavior
- **Applications:**
 - Credit Card Fraud Detection
 - Network Intrusion Detection
 - Bot detection in web traffic



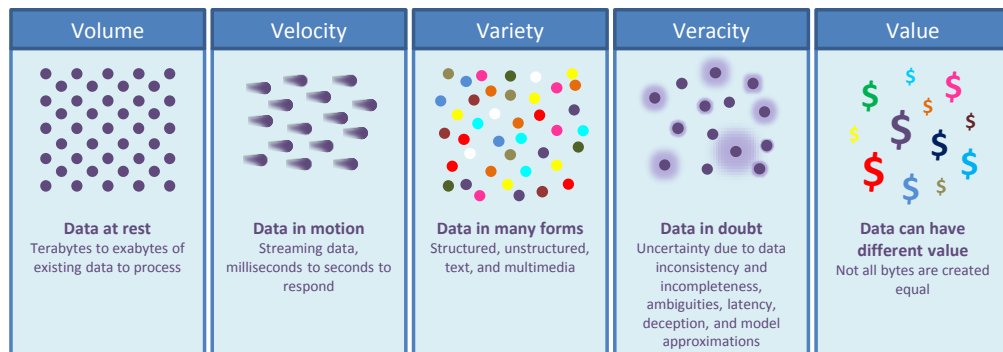
data science dojo
unleash the data scientist in you

Challenges in Data Mining



data science dojo
unleash the data scientist in you

5 Vs Of Big Data



datascience**dojo**
unleash the data scientist in you

Questions?

datascience**dojo**
unleash the data scientist in you