

INTRO TO DATA SCIENCE

LECTURE 1: DATA EXPLORATION

Arun Ahuja – aahuja11@gmail.com

Sandip Trivedi - trivedi.sandip@gmail.com

INTRO TO DATA SCIENCE

WELCOME!

COURSE MEETING:

T/TH 6:30 - 9:30

GA WEST 21ST STREET

COURSE NOTES:

[HTTP://GADATASCIENCE.COM](http://gadatasience.com)

COURSE MAILS: DAT-NYC-10@GA-GROUPS.COM

INTRO TO DATA SCIENCE

WELCOME!

I. WHAT IS DATA SCIENCE?

II. THE DATA MINING WORKFLOW

EXERCISES:

III. WORKING AT THE UNIX COMMAND LINE

INTRO TO DATA SCIENCE

I. WHAT IS DATA SCIENCE?

- A set of tools and techniques used to extract useful information from data.

- A set of tools and techniques used to extract useful information from data.
- An interdisciplinary, problem-oriented subject.

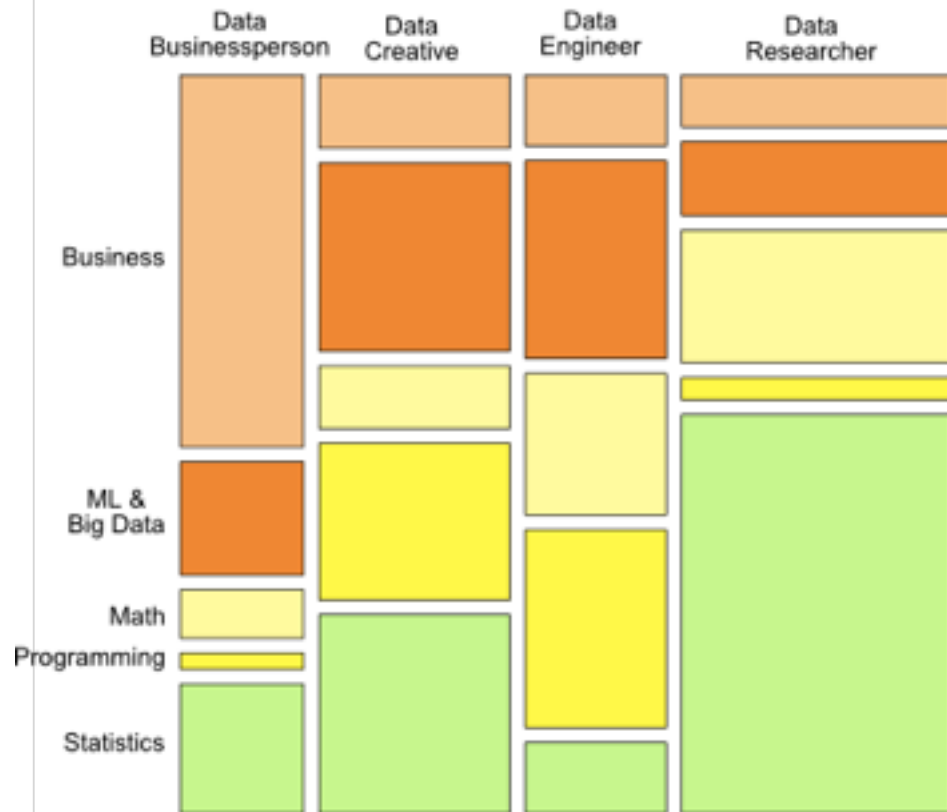
WHAT IS DATA SCIENCE?

9

Data Developer	Developer	Engineer	
Data Researcher	Researcher	Scientist	Statistician
Data Creative	Jack of All Trades	Artist	Hacker
Data Businessperson	Leader	Businessperson	Entrepreneur

WHAT IS DATA SCIENCE?

10



source: *Analyzing the Analyzers*

WHAT IS DATA SCIENCE?

11

Business	ML / Big Data	Math / OR	Programming	Statistics
Product Development	Unstructured Data	Optimization	Systems Administration	Visualization
Business	Structured Data	Math	Back End Programming	Temporal Statistics
	Machine Learning	Graphical Models	Front End Programming	Surveys and Marketing
	Big and Distributed Data	Bayesian / Monte Carlo Statistics		Spatial Statistics
		Algorithms		Science
		Simulation		Data Manipulation
				Classical Statistics



Michael E. Driscoll

@medriscoll



Following

Data scientists: better statisticians than
most programmers & better programmers
than most statisticians [@peteskomoroch](https://bit.ly/NHmRqu)



Reply



Retweet



Favorite



More



Pocket

- A set of tools and techniques used to extract useful information from data.
- An interdisciplinary, problem-solving oriented subject.
- The application of scientific techniques to practical problems.



[Your Amazon.com](#)
[Today's Deals](#)
[Gift Cards](#)
[Sell](#)
[Help](#)

[Shop by Department](#)

[Hello, Sign in to Your Account](#)
[Join Prime](#)
[Cart](#)
[List](#)

- Unlimited Instant Videos
- MP3s & Cloud Player
- 20 million songs, play anywhere
- Amazon Cloud Drive
- 5 TB of free storage
- Kindle
- Appstore for Android
- Amazon Kindle Fire
- Digital Games & Software
- Audible Audiobooks
- Books
- Movies, Music & Games
- Electronics & Computers
- Home, Garden & Tools
- Grocery, Health & Beauty
- Toys, Kids, Baby & Pet
- Clothing, Shoes & Jewelry
- Sports & Outdoors
- Automotive & Industrial
- Full Store Directory

[Instant Video](#)
[MP3 Store](#)
[Cloud Player](#)
[Kindle](#)
[Cloud Drive](#)
[Appstore for Android](#)
[Digital Games & Software](#)
[Audiobooks](#)

The Perfect Gift for Dad

Kindle Fire HD

From ~~\$199~~ **\$179**

Enter **DADDY** at checkout

Offer valid through June 6, 2013

Top Story [Outfit Trends](#) [Fast Free Shipping](#)

Amazon Fashion

Dress Shop

Our favorite wear everywhere styles from Sunday by Shelli Segal, London Tones, and more

Shop Dresses [Shop All Clothing](#)

Class of 2013

Graduation Gifts

Shop now

Try Amazon Prime

Free for 30 days

[Get Started](#)

amazon Prime

[Activate & Renew](#)

Deal of the Day

Up to 60% Off

Select BestData Memory

[See Details](#)

UP by Jawbone

Measure your activity, sleep quality, and live better

[HARLINS](#)

\$30 Off Instantly

[#30off](#)

What Other Customers Are Looking At Right Now

SanDisk Ultra 32GB MicroSDHC Memory Card

★★★★★ (2,814)

Price **\$42.99**

Samsung Galaxy Tab 2 (7-inch, Wi-Fi)

★★★★★ (2,470)

Price **\$174.00**

Kindle Fire 7" LCD Display, Wi-Fi, & Amazon Digital Navigation System

★★★★★ (8,188)

Price **\$159.00**

Kindle 7" C-In Display, Wi-Fi, & Amazon Digital Navigation System

★★★★★ (4,478)

Price **\$69.00**

Kindle Fire HD 7" Galaxy Audio, Amazon Digital Navigation System

★★★★★ (4,917)

Price **\$199.00**

Kindle Fire HD 8.9" Galaxy Audio, Amazon Digital Navigation System

★★★★★ (5,470)

Price **\$299.00**

Amazon Gift Card - Email

★★★★★ (19,180)

Price **\$50.00**

Digital Cameras Real Sellers

Nikon COOLPIX W300 16 MP CMOS

★★★★★ (1,000)

Canon PowerShot SX600 HS 16.0 MP

★★★★★ (1,000)

Canon PowerShot A1000 IS 16.0 MP

★★★★★ (1,000)

Canon PowerShot Rebel T5 18 MP CMOS

★★★★★ (719)

Canon PowerShot SX160 IS 16.1 MP

★★★★★ (1,000)

Canon PowerShot SX280 HS 12.1 MP CMOS

★★★★★ (1,000)

Canon PowerShot Rebel T6 18.0 MP CMOS

★★★★★ (1,000)

Noncredit courses and certificates

Center for Advanced Digital Applications (CADA)

There's 422 Time to...

[Your Amazon.com](#)
[Today's Deals](#)
[Gift Cards](#)
[Sell](#)
[Help](#)

[Amazon.com](#)

[Hello, Sign in](#)
[Your Account](#)
[Join Prime](#)
[Cart](#)
[List](#)

- Unlimited Instant Videos
- MP3s & Cloud Player
- 20 million songs, play anywhere
- Amazon Cloud Drive
- 5 TB of free storage
- Kindle
- Apps for Android
- Amazon Kindle
- Digital Games & Software
- Audible Audiobooks
- Books
- Movies, Music & Games
- Electronics & Computers
- Home, Garden & Tools
- Grocery, Health & Beauty
- Toys, Kids, Baby & Pet
- Clothing, Shoes & Jewelry
- Sports & Outdoors
- Automotive & Industrial
- Full Store Directory

[Instant Video](#)
[MP3 Store](#)
[Cloud Player](#)
[Kindle](#)
[Cloud Drive](#)
[Apps for Android](#)
[Digital Games & Software](#)
[Audiobooks](#)

The Perfect Gift for Dad

Kindle Fire HD

From \$199.99 to \$179.99

Enter DADS2013 at checkout

Offer valid through June 6, 2013

Class of 2013 Graduation Gifts

Try Amazon Prime free for 30 days

Get Started

Deal of the Day

Up to 60% Off

Defiant HardDisk Memory

UP by Jawbone

Measure your activity, sleep quality, and live better

FREELINE

\$30 Off Instantly

Noncredit courses and certificates

Center for Advanced Digital Applications (CADA)

There's KIDS Time to...

Top Story

Outfitting Trends

Fast Free Shipping

Amazon Fashion

Dress Shop

Our favorite wear everywhere styles from Sunday by Shelli Segal, London Tones, and more

Shop All Clothing

What Other Customers Are Looking At Right Now

<p>SanDisk Ultra SD 64GB MicroSDHC (Class 10)</p> <p>★★★★★ (4,818)</p> <p>Price \$42.99</p>	<p>Samsung Galaxy Tab 2 (7.0 inch, Wi-Fi)</p> <p>★★★★★ (2,470)</p> <p>\$7.19 to \$174.00</p>	<p>Kindle Fire 7" LCD Display, Wi-Fi, & Amazon Playful Rewards Inc.</p> <p>★★★★★ (5,188)</p> <p>\$199.00</p>	<p>Kindle 6" C Ink Display, Wi-Fi, & Amazon Playful Rewards Inc.</p> <p>★★★★★ (4,270)</p> <p>\$69.00</p>	<p>Kindle Fire HD 7" Galaxy Audio, & Amazon Playful Rewards Inc.</p> <p>★★★★★ (4,270)</p> <p>\$199.00</p>	<p>Kindle Fire HD 8.9" Galaxy Audio, & Amazon Playful Rewards Inc.</p> <p>★★★★★ (4,470)</p> <p>\$299.00</p>	<p>Amazon Gift Card - Green</p> <p>★★★★★ (19,180)</p> <p>\$50.00</p>
---------------------------------------------------------------------------------------------	----------------------------------------------------------------------------------------------	------------------------------------------------------------------------------------------------------------------	--------------------------------------------------------------------------------------------------------------	---------------------------------------------------------------------------------------------------------------	-----------------------------------------------------------------------------------------------------------------	----------------------------------------------------------------------

Digital Cameras Best Sellers

<p>Nikon COOLPIX B5000 16 MP CMOS</p> <p>★★★★★ (1,000)</p>	<p>Canon PowerShot SX600 HS 16.0 MP</p> <p>★★★★★ (1,000)</p>	<p>Canon PowerShot A1000 IS 16.0 MP</p> <p>★★★★★ (1,000)</p>	<p>Canon PowerShot T5 16 MP CMOS</p> <p>★★★★★ (1,000)</p>	<p>Canon PowerShot SX160 IS 16.1 MP</p> <p>★★★★★ (1,000)</p>	<p>Canon PowerShot SX260 HS 12.1 MP CMOS</p> <p>★★★★★ (1,000)</p>	<p>Canon PowerShot T5 16.0 MP CMOS</p> <p>★★★★★ (1,000)</p>
------------------------------------------------------------	--------------------------------------------------------------	--------------------------------------------------------------	-----------------------------------------------------------	--------------------------------------------------------------	-------------------------------------------------------------------	-------------------------------------------------------------

[Your Amazon.com](#)
[Today's Deals](#)
[Gift Cards](#)
[Sell](#)
[Help](#)

[Shop by Department](#)

[Hello, Sign in to Your Account](#)
[Join Prime](#)
[Cart](#)
[List](#)

- Unlimited Instant Videos
- MP3s & Cloud Player
- 20 million songs, play anywhere
- Amazon Cloud Drive
- 5 TB of free storage
- Kindle
- Apps for Android
- Smartphone Sync
- Free Today
- Digital Games & Software
- Audible Audiobooks
- Books
- Movies, Music & Games
- Electronics & Computers
- Home, Garden & Tools
- Grocery, Health & Beauty
- Toys, Kids, Baby & Pet
- Clothing, Shoes & Jewelry
- Sports & Outdoors
- Automotive & Industrial
- Full Store Directory

[Instant Video](#)
[MP3 Store](#)
[Cloud Player](#)
[Kindle](#)
[Cloud Drive](#)
[Apps for Android](#)
[Digital Games & Software](#)
[Audiobooks](#)

The Perfect Gift for Dad

Kindle Fire HD

From \$199 From \$179

Enter DADS179 at checkout

Offer valid through June 6, 2013

Class of 2013 Graduation Gifts

Shop now

Try Amazon Prime free for 30 days

Get Started

Deal of the Day

Up to 60% Off

Defiant Backpack Memory

Shop now

UP by Jawbone

Measure your activity, sleep quality, and live better

FREE

Top Story

Outfit Trends

Fast Free Shipping

Amazon Fashion

Dress Shop

Our favorite wear everywhere styles from Sunday by Shelli Segal, London Tones, and more

Shop Now

What Other Customers Are Looking At Right Now

SanDisk Ultra 32GB microSDHC Card

4.9 (1,111)

\$42.99

Samsung Galaxy Tab 2 (7-inch, Wi-Fi)

4.3 (1,111)

\$179.00

Kindle Fire 7" LCD Display, Wi-Fi, & Amazon Digital Navigation bar

4.8 (1,111)

\$159.00

Kindle 7" C Ink Display, Wi-Fi, & Amazon Digital Navigation bar

4.8 (1,111)

\$69.00

Kindle Fire HD 7", Galaxy Audio, & Amazon Digital Navigation bar

4.8 (1,111)

\$199.00

Kindle Fire HD 8.9", Galaxy Audio, & Amazon Digital Navigation bar

4.8 (1,111)

\$289.00

Amazon Gift Card - Credit

4.9 (1,111)

\$50.00

Digital Cameras Best Sellers

Nikon COOLPIX W300 16 MP CMOS

4.9 (1,111)

Canon PowerShot SX600 HS 16.0 MP

4.9 (1,111)

Canon PowerShot A2000 HS 16.0 MP

4.9 (1,111)

Canon PDS Rebel T5 18 MP CMOS

4.9 (1,111)

Canon PowerShot SX700 HS 16.1 MP

4.9 (1,111)

Canon PowerShot SX600 HS 16.1 MP CMOS

4.9 (1,111)

Canon PDS Rebel T5 18.0 MP CMOS

4.9 (1,111)

Noncredit courses and certificates

Center for Advanced Digital Applications (CADA)

There's Still Time to...



Roll over image to zoom in



[See 1 customer image](#)

[Share your own customer images](#)

Star Trek [Blu-ray] (2009)

[Chris Pine](#) (Actor), [Zachary Quinto](#) (Actor), [J.J. Abrams](#) (Director) | Rated: PG-13 | Format: Blu-ray

[View reviews](#) (2,040 customer reviews)

List Price: **\$22.68**

Price: **\$9.99** & **FREE Shipping** on orders over \$25. [Details](#)

You Save: **\$12.69** (57%)

In Stock.

Ships from and sold by [Amazon.com](#). Gift-wrap available.

Want it Wednesday, June 5? Order within 30 hrs 11 mins and choose **One-Day Shipping** at checkout. [Details](#)

23 new from \$9.98 **19 used** from \$9.78 **1 collectible** from \$49.99

Watch Instantly with amazon instant video		Rent	Buy		
Star Trek (2009)		\$2.99	\$9.99		
Other Formats & Versions		Amazon Price	New from	Used from	
	Blu-ray	1-Disc Version	\$9.99	\$9.98	\$9.78
	DVD	Single-Disc Edition	\$8.49	\$5.65	\$3.29



This week only, save up to 62% on [Farscape: The Complete Series](#) in our Deal of the Week. Offer ends June 8, 2013. [Learn more](#)

Frequently Bought Together



+



+



Price for all three: **\$66.47**

[Add all three to Cart](#)

[Add all three to Wish List](#)

Some of these items ship sooner than the others. [Show details](#)

- ☒ **This item:** Star Trek [Blu-ray] ~ Chris Pine Blu-ray **\$9.99**
- ☒ Star Trek Into Darkness (Blu-ray 3D + Blu-ray + DVD + Digital Copy) ~ Chris Pine Blu-ray **\$24.99**
- ☒ Iron Man 3 (Two-Disc Blu-ray / DVD + Digital Copy) ~ Robert Downey Jr. Blu-ray **\$31.49**

What Other Items Do Customers Buy After Viewing This Item?



Star Trek Into Darkness (Blu-ray 3D + Blu-ray + DVD + Digital Copy) ~ Chris Pine Blu-ray
[View reviews](#) (199)
\$24.99



Star Trek: Original Motion Picture Collection (Star Trek I, II, III, IV, V, VI + The Captain's Summit Bonus Disc) [Blu-ray] ~ William Shatner Blu-ray
[View reviews](#) (371)
\$53.56



Sin City (Two-Disc Theatrical & Recut, Extended, and Unrated Versions) [Blu-ray] ~ Jessica Alba Blu-ray
[View reviews](#) (933)
\$4.99

Quantity:

Yes, I want **FREE Two-Day Shipping** with [Amazon Prime](#)

[Add to Cart](#)

or

[Sign in](#) to turn on 1-Click ordering.

[Add to Wish List](#)

Sell Us Your Item

For up to a **\$2.60** Gift Card

[Trade in](#)

[Learn more](#)

More Buying Choices

[TechShowable](#) [Add to Cart](#)
\$19.99 & **FREE Shipping** on orders over \$25. [Details](#)

43 used & new from \$9.78

Have one to sell? [Sell on Amazon](#)

[Share](#) [Facebook](#) [Twitter](#) [Pinterest](#)



Roll over image to zoom in



See 1 customer image

Share your own customer images

Star Trek [Blu-ray] (2009)

Chris Pine (Actor), Zachary Quinto (Actor), J.J. Abrams (Director) | Rated: PG-13 | Format: Blu-ray

4.4 (2,040 customer reviews)

List Price: ~~\$22.68~~

Price: \$9.99 & **FREE Shipping** on orders over \$25. [Details](#)

You Save: \$12.69 (57%)

In Stock.

Ships from and sold by Amazon.com. Gift-wrap available.

Want it Wednesday, June 5? Order within 30 hrs 11 mins and choose **One-Day Shipping** at checkout. [Details](#)

23 new from \$9.98 **19 used** from \$9.78 **1 collectible** from \$49.99

Watch instantly with amazon instant video		Rent	Buy
Star Trek (2009)		\$2.99	\$9.99
Other Formats & Versions		Amazon Price	New from Used from
① Blu-ray	1-Disc Version	\$9.99	\$9.98 \$9.78
② DVD	Single-Disc Edition	\$8.49	\$5.65 \$3.29



This week only, save up to 62% on [Farscape: The Complete Series](#) in our Deal of the Week. Offer ends June 8, 2013. [Learn more](#)

Frequently Bought Together



+



+



Price for all three: **\$66.47**

[Add all three to Cart](#)

[Add all three to Wish List](#)

Some of these items ship sooner than the others. [Show details](#)

- ✓ **This item:** Star Trek [Blu-ray] ~ Chris Pine Blu-ray \$9.99
- ✓ Star Trek Into Darkness (Blu-ray 3D + Blu-ray + DVD + Digital Copy) ~ Chris Pine Blu-ray \$24.99
- ✓ Iron Man 3 (Two-Disc Blu-ray / DVD + Digital Copy) ~ Robert Downey Jr. Blu-ray \$31.49

What Other Items Do Customers Buy After Viewing This Item?



Star Trek Into Darkness (Blu-ray 3D + Blu-ray + DVD + Digital Copy) ~ Chris Pine Blu-ray
★★★★☆ (199)
\$24.99



Star Trek: Original Motion Picture Collection (Star Trek I, II, III, IV, V, VI + The Captain's Summit Bonus Disc) [Blu-ray] ~ William Shatner Blu-ray
★★★★☆ (571)
\$53.56



Sin City (Two-Disc Theatrical & Recut, Extended, and Unrated Versions) [Blu-ray] ~ Jessica Alba Blu-ray
★★★★☆ (933)
\$4.99

Quantity:

Yes, I want **FREE Two-Day Shipping** with [Amazon Prime](#)

[Add to Cart](#)

or

[Sign in](#) to turn on 1-Click ordering.

[Add to Wish List](#)

Sell Us Your Item

For up to a **\$2.60 Gift Card**

[Trade in](#)

[Learn more](#)

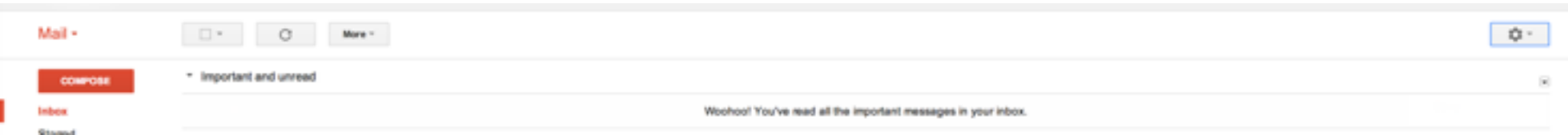
More Buying Choices

Technoshade [Add to Cart](#)
\$19.99 & **FREE Shipping** on orders over \$25. [Details](#)

43 used & new from \$9.78

Have one to sell? [Sell on Amazon](#)

Share [Facebook](#) [Twitter](#) [Pinterest](#) [Google+](#)



II. THE DATA SCIENCE WORKFLOW

acquire parse filter mine represent refine interact













‣ From Jeff Hammerbacher:

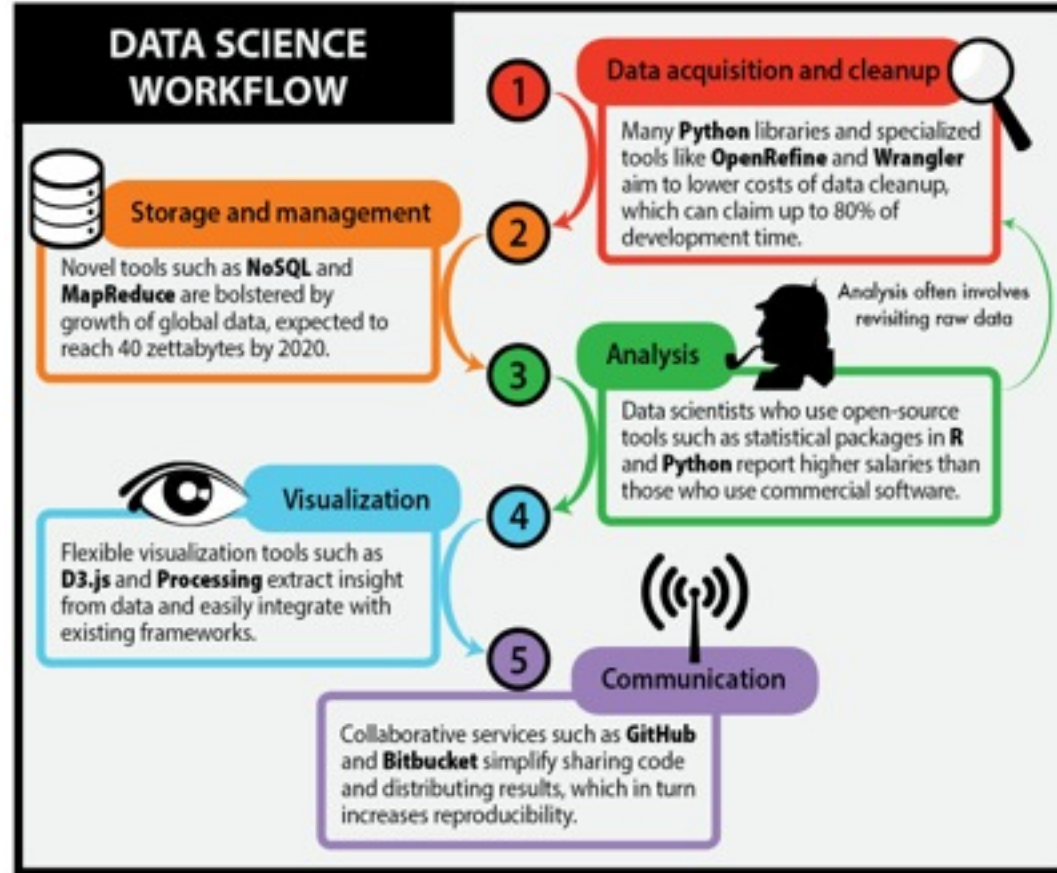
1. Identify problem
2. Instrument data sources
3. Collect data
4. Prepare data (integrate, transform, clean, impute, filter, aggregate)
5. Build model
6. Evaluate model
7. Communicate results

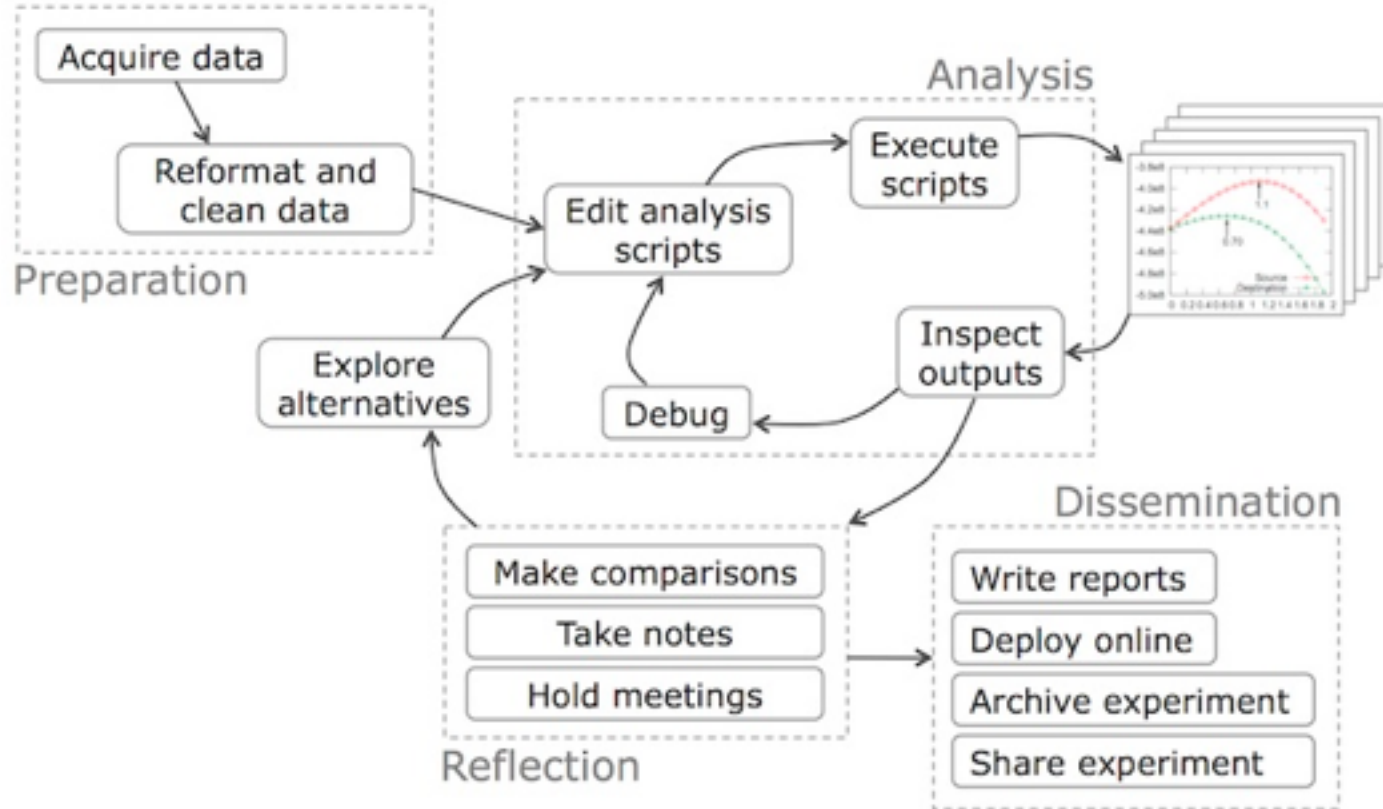
‣ From Jeff Hammerbacher:

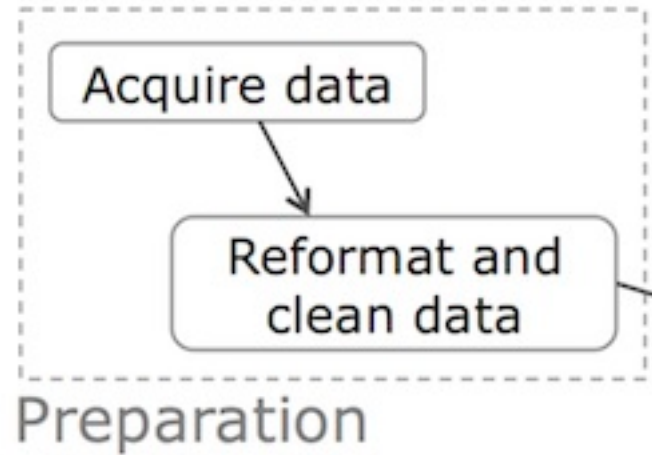
1. Identify problem
2. Instrument data sources
3. Collect data
4. Prepare data (integrate, transform, clean, impute, filter, aggregate)
5. Build model
6. Evaluate model
7. Communicate results

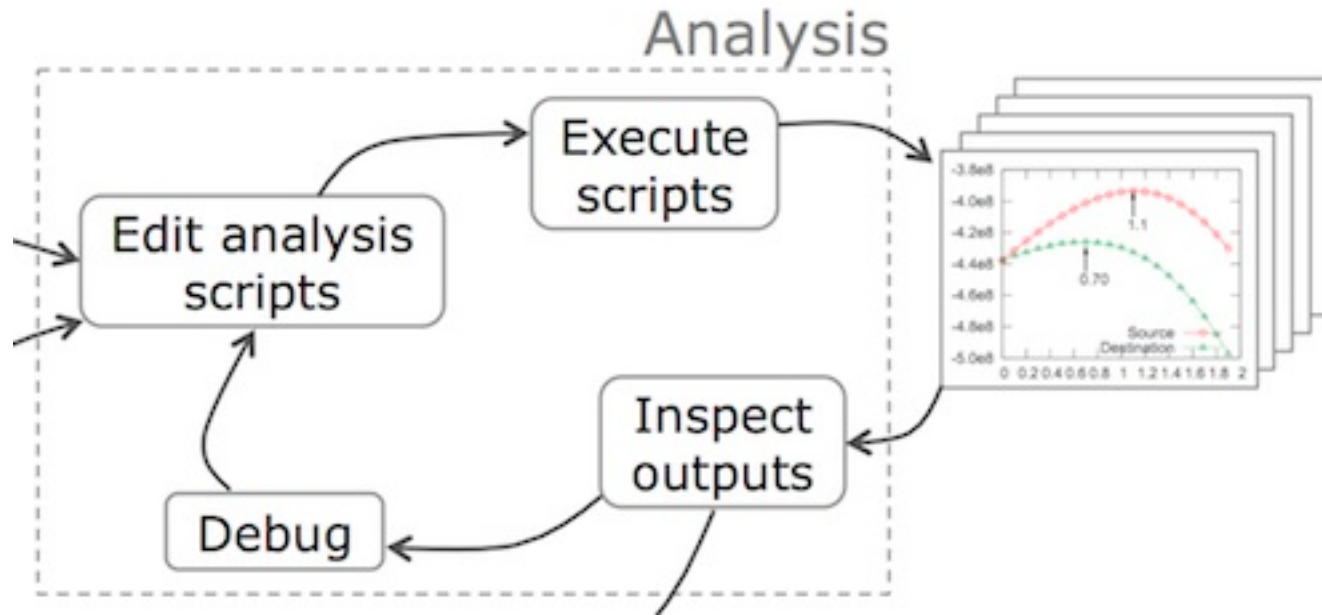
From Dataists Blog

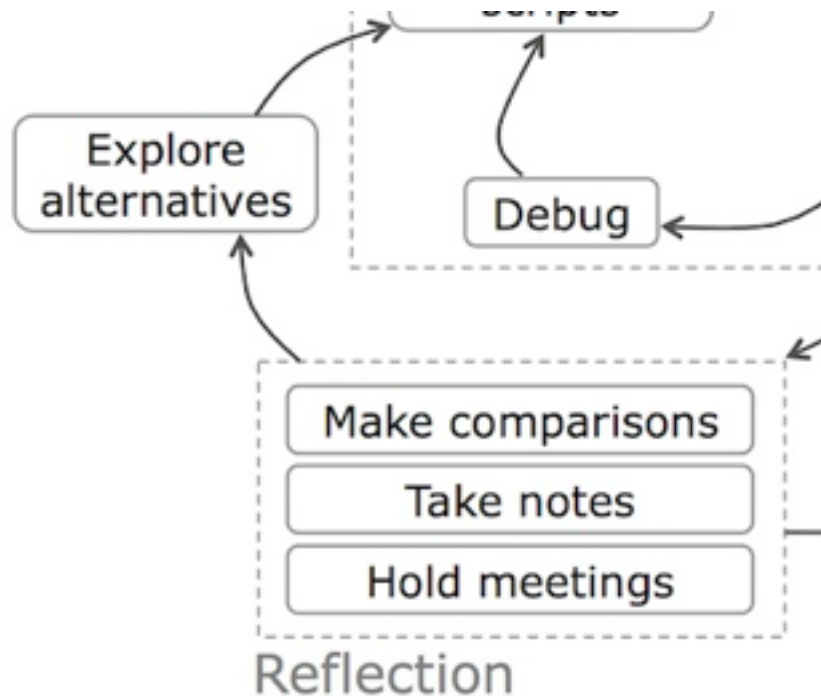
1. Obtain
2. Scrub
3. Explore
4. Model
5. Interpret











PROBLEM: WHAT ARE THE LEADING INDICATORS THAT A USER WILL MAKE A NEW PURCHASE?

PROBLEM: WHAT ARE THE LEADING INDICATORS THAT A USER WILL MAKE A NEW PURCHASE?

1. Collect data around user retention, user actions within the product, potentially find data outside of company

PROBLEM: WHAT ARE THE LEADING INDICATORS THAT A USER WILL MAKE A NEW PURCHASE?

1. Collect data around user retention, user actions within the product, potentially find data outside of company
2. Extract aggregated values from raw data
 1. How many times did a user share through Facebook within a week? A month?
 2. How often did they open up our emails?

PROBLEM: WHAT ARE THE LEADING INDICATORS THAT A USER WILL MAKE A NEW PURCHASE?

1. Collect data around user retention, user actions within the product, potentially find data outside of company
2. Extract aggregated values from raw data
 1. How many times did a user share through Facebook within a week? A month?
 2. How often did they open up our emails?
3. Examine data to find common distributions and correlations

PROBLEM: WHAT ARE THE LEADING INDICATORS THAT A USER WILL MAKE A NEW PURCHASE?

1. Collect data around user retention, user actions within the product, potentially find data outside of company
2. Extract aggregated values from raw data
 1. How many times did a user share through Facebook within a week? A month?
 2. How often did they open up our emails?
3. Examine data to find common distributions and correlations
4. Extract new meaning to predict if a user would purchase again or not

PROBLEM: WHAT ARE THE LEADING INDICATORS THAT A USER WILL MAKE A NEW PURCHASE?

1. Collect data around user retention, user actions within the product, potentially find data outside of company
2. Extract aggregated values from raw data
 1. How many times did a user share through Facebook within a week? A month?
 2. How often did they open up our emails?
3. Examine data to find common distributions and correlations
4. Extract new meaning to predict if user would purchase again
5. Share results (and probably also go back to the drawing board)

Let's build an analytics team:

You are on of the following organizations:

GA, Pinterest, Tesla, Pixar, Blue Bottle Coffee

Answer the following questions:

1. Define the top priorities of the organization
2. What products could you build or studies could you run?
3. What data would you need to collect?
4. What will your top challenges be?

III. WORKING AT THE UNIX COMMAND LINE

LET'S TAKE A LOOK AT THE 538 DATASET

INTRO TO DATA SCIENCE

DISCUSSION