# Technology Fundamentals for Analytics

Jason Kuruzovich

# Overview

- Introductions
- Let's Get Excited about Data Science
- What do Data Scientists Do?
- Course Syllabus
- Lab 1. Development Environment

Next Startup Tech Valley Event is September 2 at 5:30 at Brown's Revolution Hall
See more at www.startuptechvalley.org

# Me

- Director of the Severino Center for Technological Entrepreneurship

- Associate Professor of Business Analytics

- Research on online markets and entrepreneurship

# You

- What is your background?
- What are you looking for out of this class/program?
- What type of job do you want?
- What are your hobbies?

There have been profound changes in technology and the information processes define our society
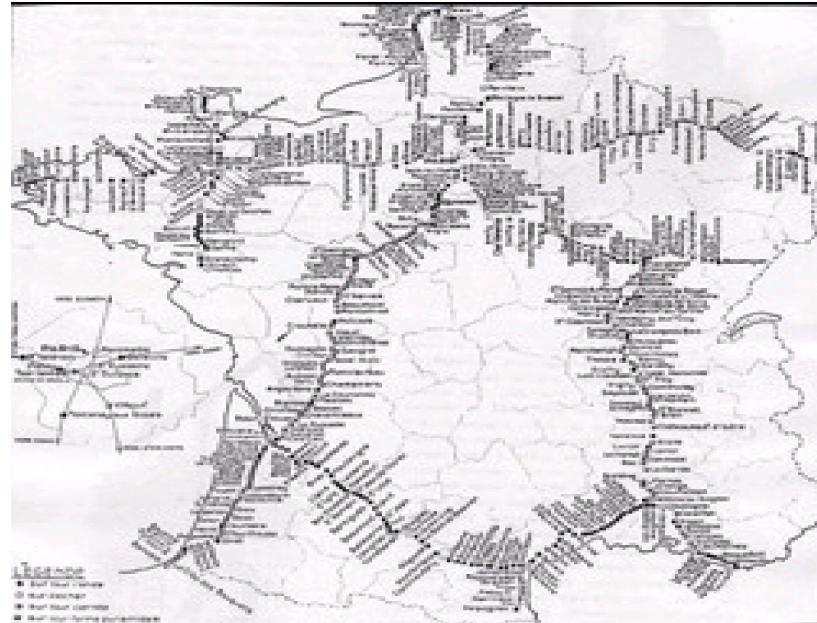
# Internet 0.1 beta (18<sup>th</sup> Century)

# Internet 0.1 Beta (18<sup>th</sup> Century)

Chain of towers or optical telegraph capable of transmitting 1-3 symbols per min

Towers 5-20 KM apart

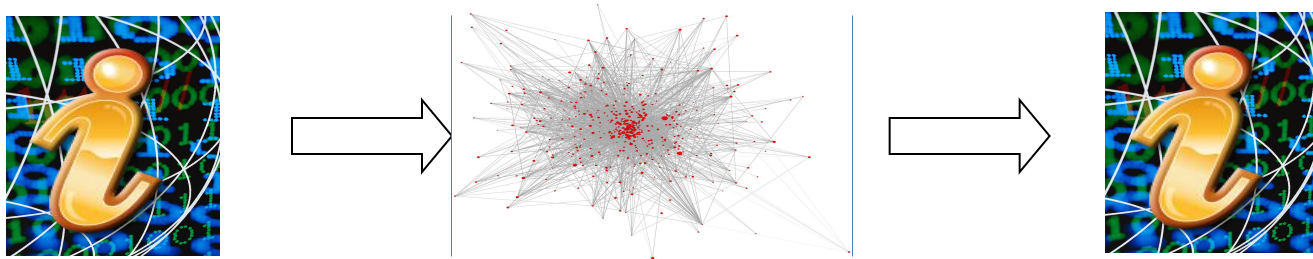"We create as much information in two days now as we did from the dawn of man through 2003."

-Eric Schmidt, Former CEO of Google

# Information Economy

TRADITIONAL PRODUCTION PROCESS
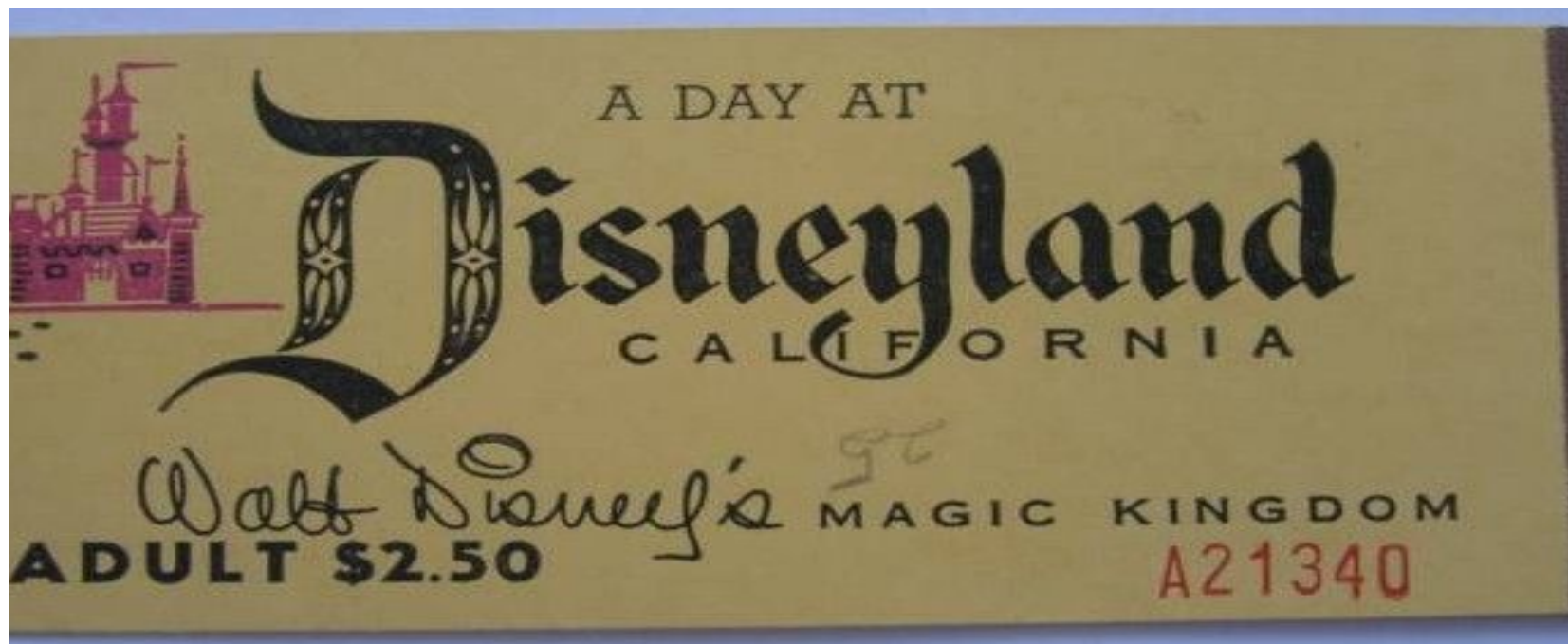


INFORMATION BASED BUSINESS PROCESS



INFORMATION TECHNOLOGY

# What is driving the growth of information and data?
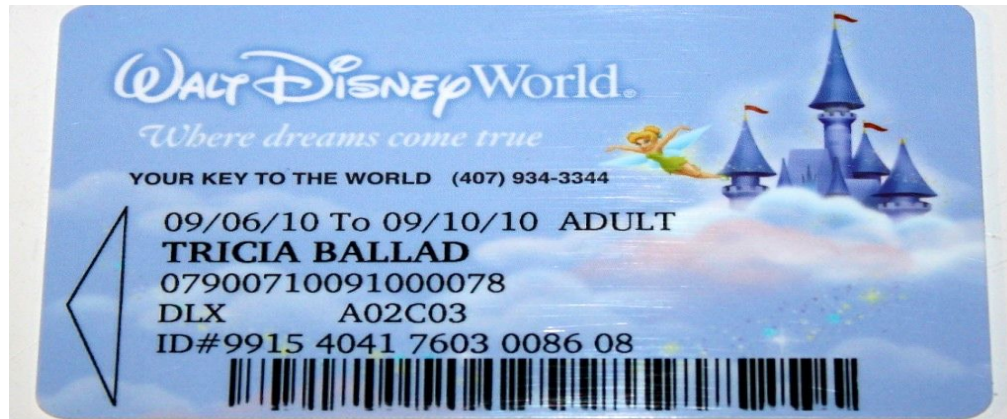
Consider the evolution of a company like Disney…

# Disney

**ROLE OF DATA: How many tickets did we sell?**

# Disney – Data Warehouse Stage

**ROLE OF DATA: How much did our customers spend? How can we understand different customer types?**

# Disney – Big Data

**ROLE OF DATA:** **What path did customers take through the park, when did they leave? How long did they stand in line? When did they spend money on souvenirs and where? How often did they go to the bathroom and did they have to wait? How long did they spend at dinner in the Mexican pavilion compared with the German pavilion? How does the speed of entry correlate with tipping behavior?**

# Electronic Commerce

- TRADITIONAL SHOPPING
- TRAFFIC AND POS DATA



- EVERY CLICK
- HISTORICAL PURCHASE
- REMARKETING

# What if the online and offline merge?

https://youtu.be/uiDMlFycNrw

# Cost per Megabase of DNA Sequence

Moore's Law

National Human Genome Research Institute
genome.gov/sequencingcosts

# Big Data and Astronomy



The Murchison Widefield Array is the first Square Kilometre Array precursor to enter full operations, generating a vast torrent of information that needs to be stored for later retrieval by researchers.

http://www.skatelescope.org/news/pawsey-centre/

"To store the Big Data the MWA produces, you'd need almost three 1 TB hard drives every two hours"

# Social Media

# The Facebook Social Network

# Facebook

- 300 Petabyte Data warehouse

- 500+ Terabytes new data each day

# Online Video

# YouTube

- 60 hours of video per minute

- 4 billion views per day

- 800 million unique users

# Search

# Google Flu Trends



How Google Flu Trends Works

CDC ILI Data     Query 3 ✓

http://www.google.org/flutrends/about/how.html

# Internet of Things

# Technology roadmap: The Internet of Things



Technology Reach (y-axis)

Time (x-axis): 2000, 2010, 2020

Software agents and advanced sensor fusion

Miniaturisation, power-efficient electronics, and available spectrum

Teleoperation and telepresence: Ability to monitor and control distant objects

Physical-World Web

Ability of devices located indoors to receive geological signals

Locating people and everyday objects

Ubiquitous Positioning

Cost reduction leading to diffusion into 2nd wave of applications

Surveillance, security, healthcare, transport, food safety, document management

Vertical-Market Applications

Demand for expedited logistics

RFID tags for facilitating routing, inventorying, and loss prevention

Supply-Chain Helpers

Source: SRI Consulting Business Intelligence

Each day GE gathers "50 million pieces of data from 10 million sensors, off equipment worth $1 trillion."

| Product Company | → | Data Services Company |

Almost the entire corpus of literature is now digital…

# The Expression of Emotions in 20th Century Books



"using the data set provided by Google that includes word frequencies in roughly 4% of all books published up to the year 2008. We find evidence for distinct historical periods of positive and negative
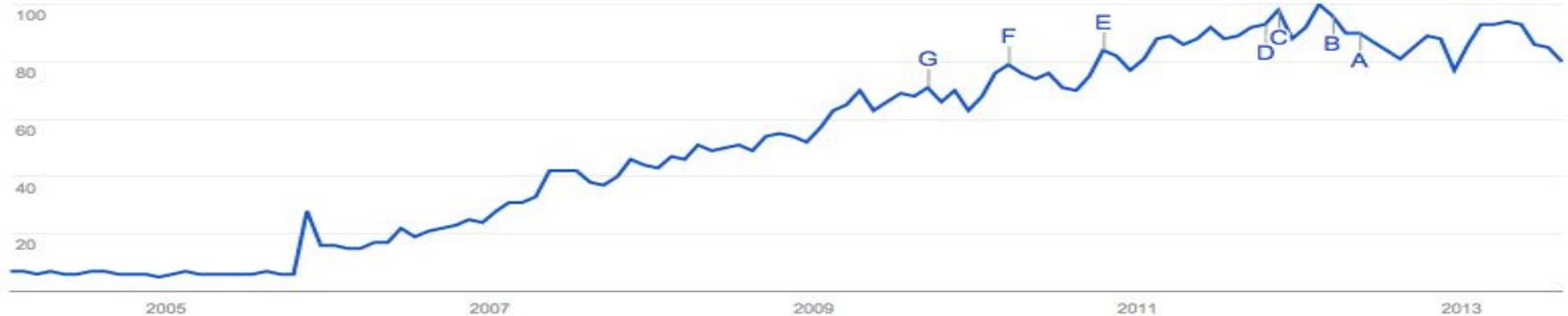
# What do we mean by "Analytics"?
## (Term Frequency in News)



Source: http://www.google.com/trends/explore?q=analytics#q=analytics&cmpt=q

# What do we mean by "Analytics"?

# Analytics as Data Apps

- The web is full of "data-driven apps."

*"The thread that ties most of these applications together is that data collected from users provides added value. Whether that data is search terms, voice samples, or product reviews, the users are in a feedback loop in which they contribute to the products they use. That's the beginning of data science."*

Prediction is a Great
Business Model!
Even if it doesn't
use analytics!!!

# Predict Who is Going to Be a Good Hire

# Predict the Best Type/Time of Post

# Predict what Terms Can Help Get You Hired (RPI STARTUP!!)

Track Company Growth
http://mattermark.com



Find which Customers Will Churn
http://framed.io



A/B Testing for Mobile Apps
https://taplytics.com



Segment Customers
http://www.segmentify.com

Predict the players who are the best value…

# Predict Where to Advertise

- "The 2012 campaign took advantage of advanced set-top-box monitoring technology to figure out what shows the voters they wanted to reach were watching and when, resulting in a smarter and cheaper — if potentially more invasive — way to beam commercials into their homes. The system gave Obama a significant advantage over Mitt Romney, according to Democrats and many Republicans (at least those who were not on Romney's media team)."

Source:
http://www.nytimes.com/2013/06/23/magazine/the-obama-campaigns-digital-masterminds-cash-in.html?pagewanted=all&_r=0

"**Analytics** is the discovery and communication of meaningful patterns in data."

-Wikipedia

The goal of this course will be to provide the <span style="color:red">technical foundation</span> to enable students to become <span style="color:red">data scientists.</span>

# What do "Data Scientists" do?



Mentions of Data Scientist in Google from Google Trends

# Data Scientist

- "A data scientist is someone who can obtain, scrub, explore, model and interpret data, blending hacking, statistics and machine learning. Data scientists not only are adept at working with data, but appreciate data itself as a first-class product."

  - Hilary Mason, chief scientist at bit.ly

# Data Scientist

Data science requires skills ranging from traditional computer science to mathematics to art. Describing the data science group he put together at Facebook (possibly the first data science group at a consumer-oriented web property), Jeff Hammerbacher said:

*"… on any given day, a team member could author a multistage processing pipeline in Python, design a hypothesis test, perform a regression analysis over data samples with R, design and implement an algorithm for some data-intensive product or service in Hadoop, or communicate the results of our analyses to other members of the organization."*

"**Analytics** is the discovery and communication of meaningful patterns in data."

-Wikipedia

Data scientists do **analytics**.

# Well…how do you "do analytics"?

# THE CRISP-DM PROCESS MODEL



Cross Industry Standard Process for Data Mining
https://en.wikipedia.org/wiki/Cross_Industry_Standard_Process_for_Data_Mining

# Business Understanding

- Project objectives and requirements
  - Are we trying to *predict* or *understand*?
  - Do we want to reduce churn, segment customers, A/B test designs…etc.
- Converting knowledge into a data mining problem definition and a preliminary plan
  - What data is needed/accessible?

# Data Understanding

- Data properties, type, distribution,
  - Understand data

- Data quality
  - Are there issues with the data? Is it as we would expect?

Visualization can be your friend for data understanding

# Data Preparation

- Clean data
  - Missing values? Outliers?

- Match data from different sources
  - Common key? If not can we match text fields.

- Feature creation
  - Does data need to be processed to enable new insights?

- Aggregate, sample and subset
  - Do we want to do analysis on entire dataset, an aggregation?

# Modeling

- Feature selection
  - Which of the identified features should be included in the model
- Splitting (fold) data
  - How should data be split to train the algorithm
- Algorithm selection and tuning
  - What algorithm should be used and what parameters set

# Evaluation

- Statistical output (Understanding)
  - What data are meaningful in predicting the variable of interest
- Prediction
  - How well does the data predict the desired outcome in the training and test datasets

# Deployment

- Implement prediction in a business process or in an application
  - Product recommendation
  - Job applicant
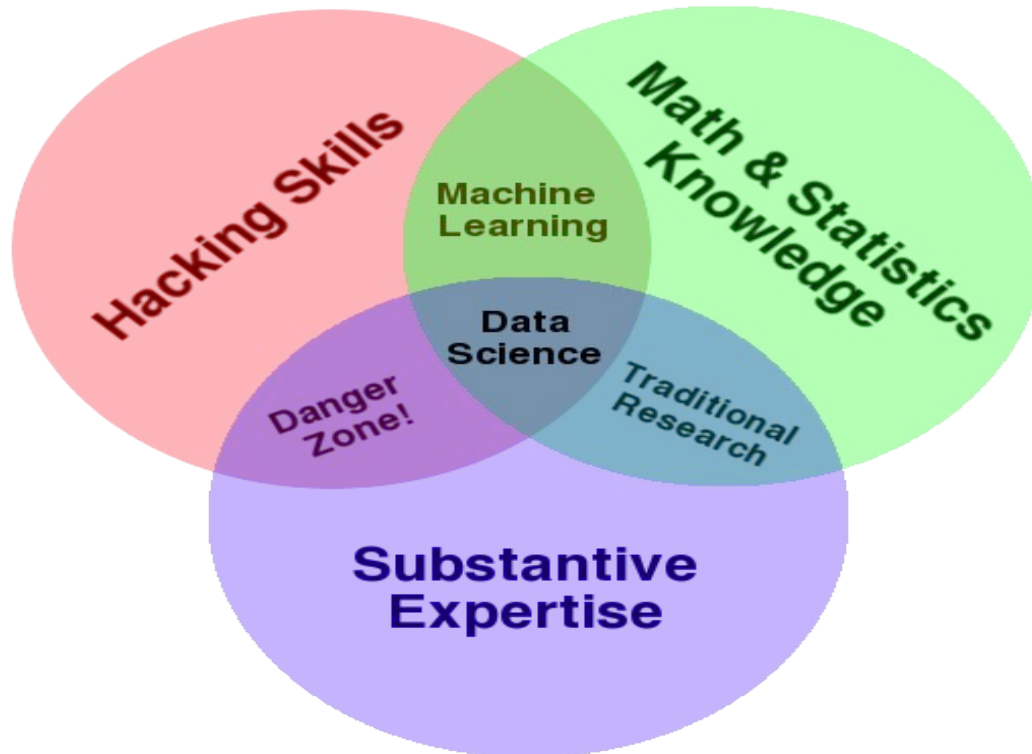  - Customer segmentation

# Data Analytics in the Wild

In a small group of 3-4, quickly discuss what you select for the data analytics in the wild assignment

Then go through the process (at a high level) for business understanding (and even a bit of data understanding) for the most interesting one in the group

Be prepared to report back to the class (~1 minute)

# What skills are needed as a data scientist?

# Data Science Venn Diagram

# Key Tools of the Data Scientist

- Data Munging - parsing, scraping, and formatting data

- Statistics - traditional analysis you're used to thinking about

- Visualization - graphs, tools, etc.

# 8 Skills to Get You Hired as a Data Scientist

1. Basic Tools
2. Basic Statistics
3. Machine Learning
4. Multivariable Calculus and Linear Algebra
5. Data Munging
6. Data Visualization & Communication
7. Software Engineering
8. Thinking Like A Data Scientist

http://blog.udacity.com/2014/11/data-science-job-skills.html

# Abstractions vs Tools

- *Abstractions* of data science
  - Matrices and linear algebra
  - Relations and relational algebra
  - MapReduce
  - Feature selection in Visualization
- Tools
  - Python
  - R
  - SAS
  - SQL/MySQL
  - Hadoop (MapReduce)
  - Tableau (Visualization)

# Syllabus and Lab