

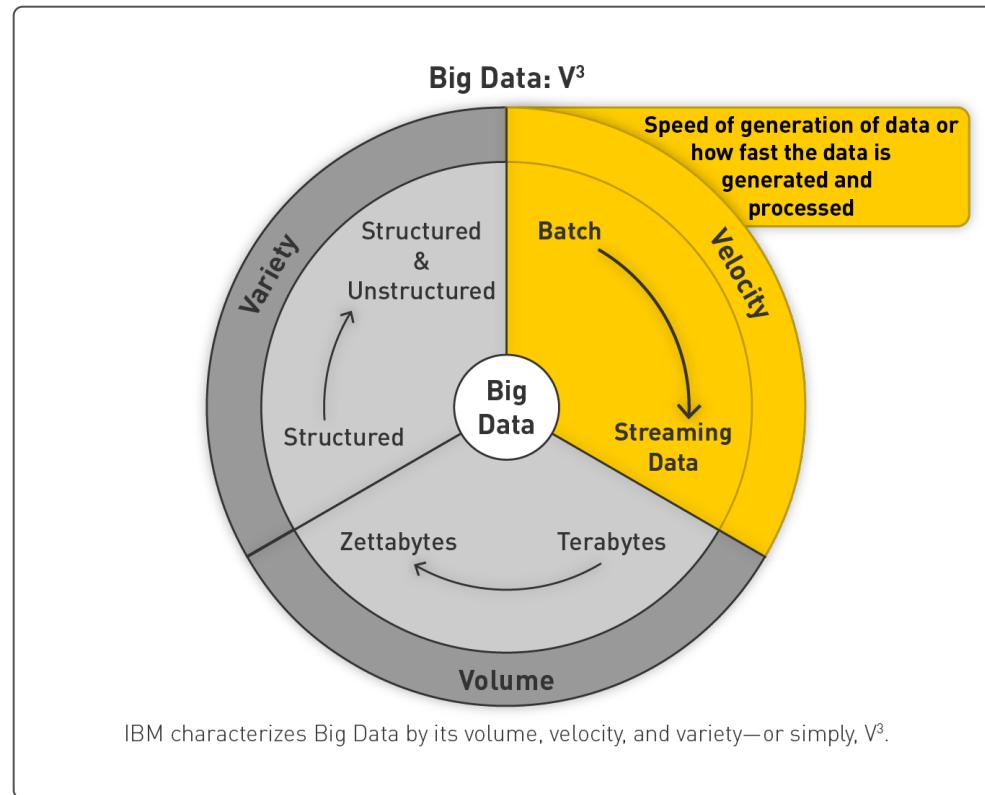
Big Data Definition

Big data is broad term for data sets so large or complex that traditional data-processing applications are **inadequate**:

- Processing
 - Laptop (8–16 GB of memory, 1 TB hard drive) overwhelmed with 4–5 GB
- Storage
 - Laptop only 1 TB
- Throughput
 - Three hours to read 1 TB on laptop

Other Challenges

- Analysis
- Capture
- Data curation
- Search
- Sharing
- Storage
- Transfer
- Visualization
- Security
- Information privacy



Sources Driving Big Data

It's All Happening Online

Every:
Click
Ad impression
Billing event
Fast forward, pause, . . .
Friend request
Transaction
Network message
Fault
...



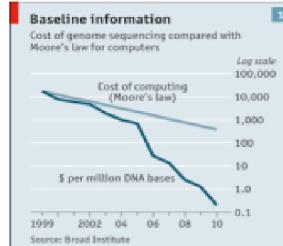
User Generated (Web, Social & Mobile) Quantified Self

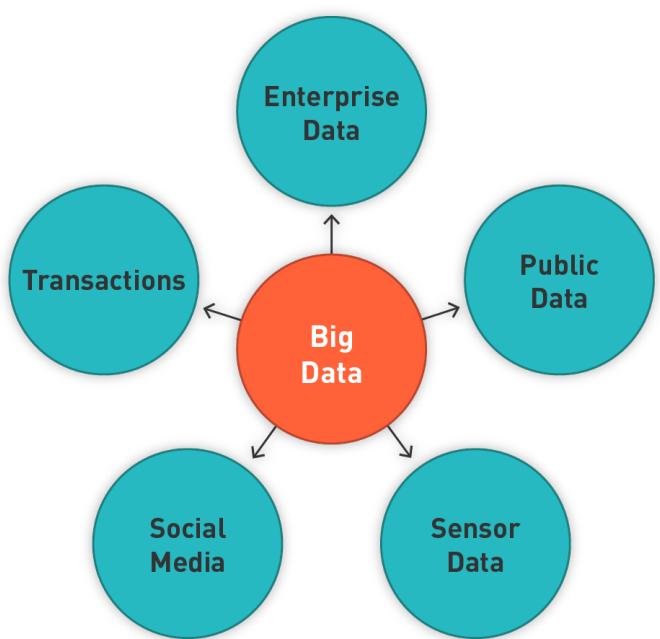


Internet of Things/M2M



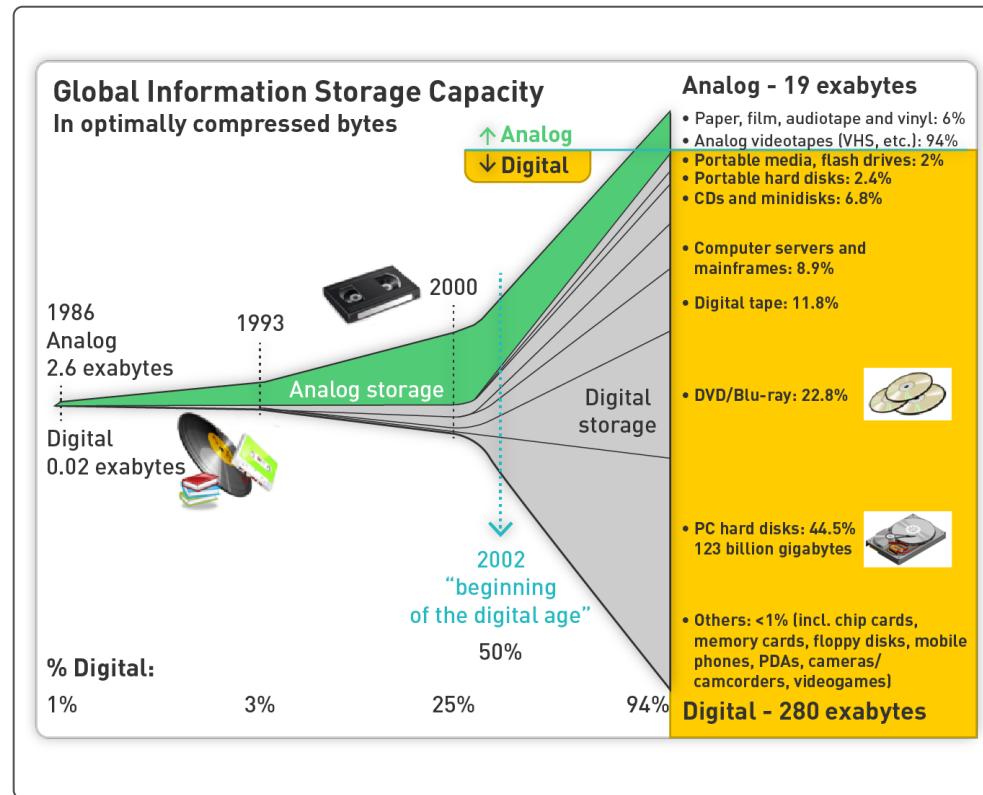
Scientific Computing





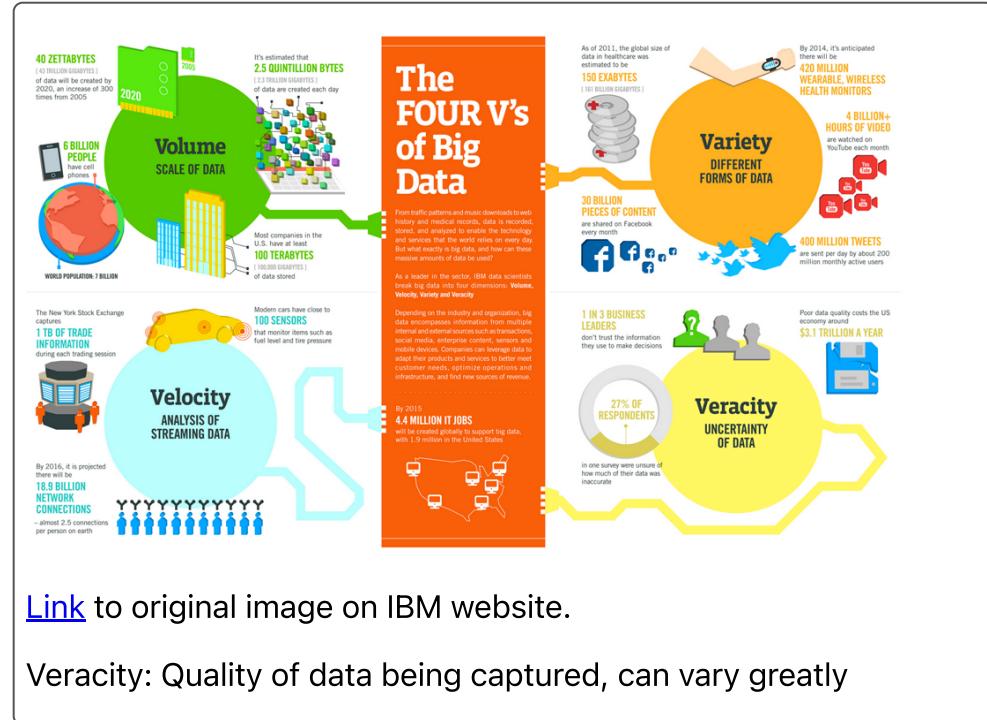
Areas of Application

New York Stock Exchange	1 TB of trade data daily
Facebook	Over 250 billion images, 250 PB of data
Ancestry.com	2.5 PB of data
Internet archives	2 PB of data
Hadron collider	15 PB of data yearly



Data in Zettabytes (ZB)

- 1 ZB = 10^{21} = 1 billion laptops (one per person)
- In 2005: 0.12 ZB
- In 2007: 0.28 ZB
- In 2020: 40 ZB
- In context, today everyone's hard drive with 1 TB
 - Data storage today: 1–2 TB data per person (1–2 billion hard drives)
 - By 2020: 40 TB per person



Applications of Big Data

- Society based (human-centric)
- Internet of things (machine-centric)

Societal and Personal Sources of Data

- Social, professional
- Quantified self (eating, sleeping, exercising)
- Voting
- Health care

Online Society

- 3 billion people online today
- 1.5 billion on a social network
 - Facebook
 - Twitter
 - LinkedIn
 - WeChat
- Sharing huge volume of data



140 billion
images, 6 billion
added monthly



6 billion images

1 billion images
served daily

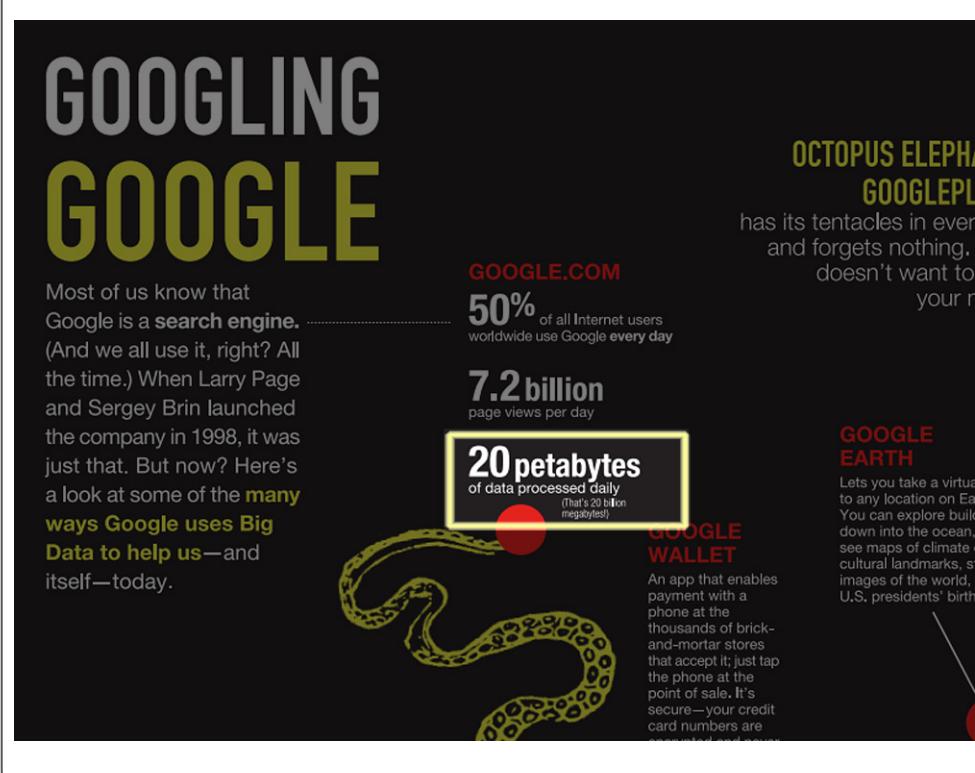


72 hours
uploaded
every minute

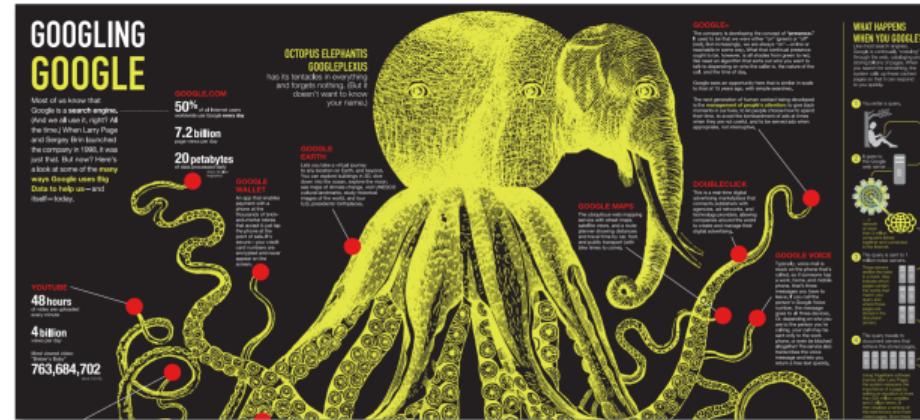


3.5 trillion
photographs

90% of net traffic will be visual.



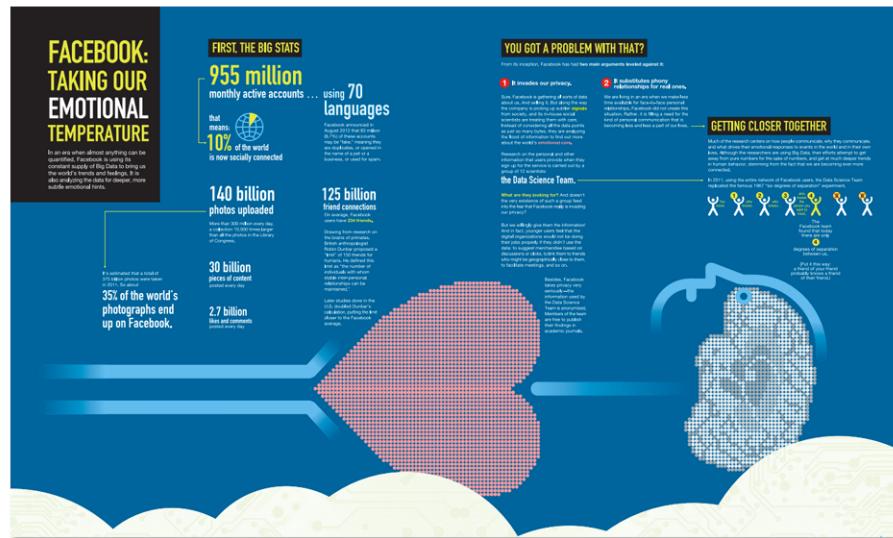




Data from 2012

Facebook Big Data Stats

- 2.5 billion content items shared daily
- 1.5 billion people on network
- Over a trillion edges in graph
- 2.7 billion "likes" daily
- 105 TB data scanned via Hive every 30 minutes



the
is
ves. →

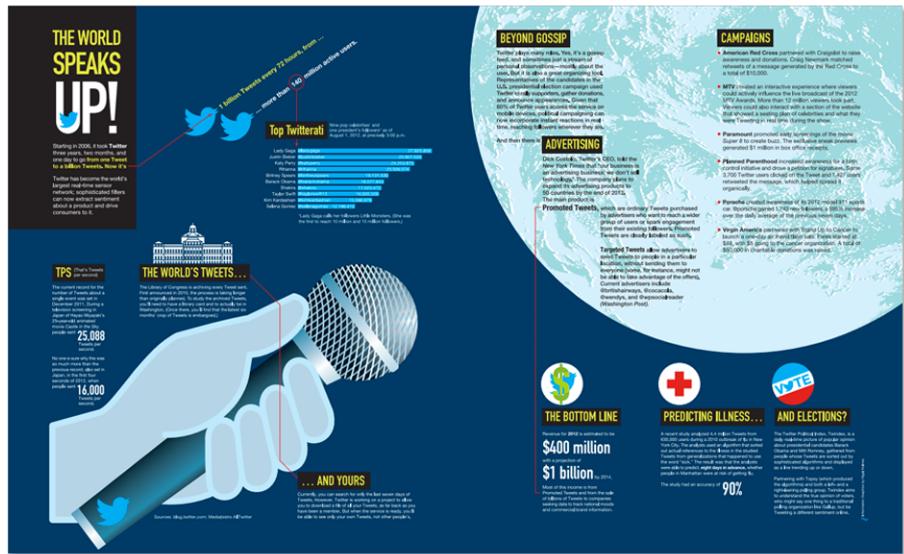
GETTING CLOSER TOGETHER

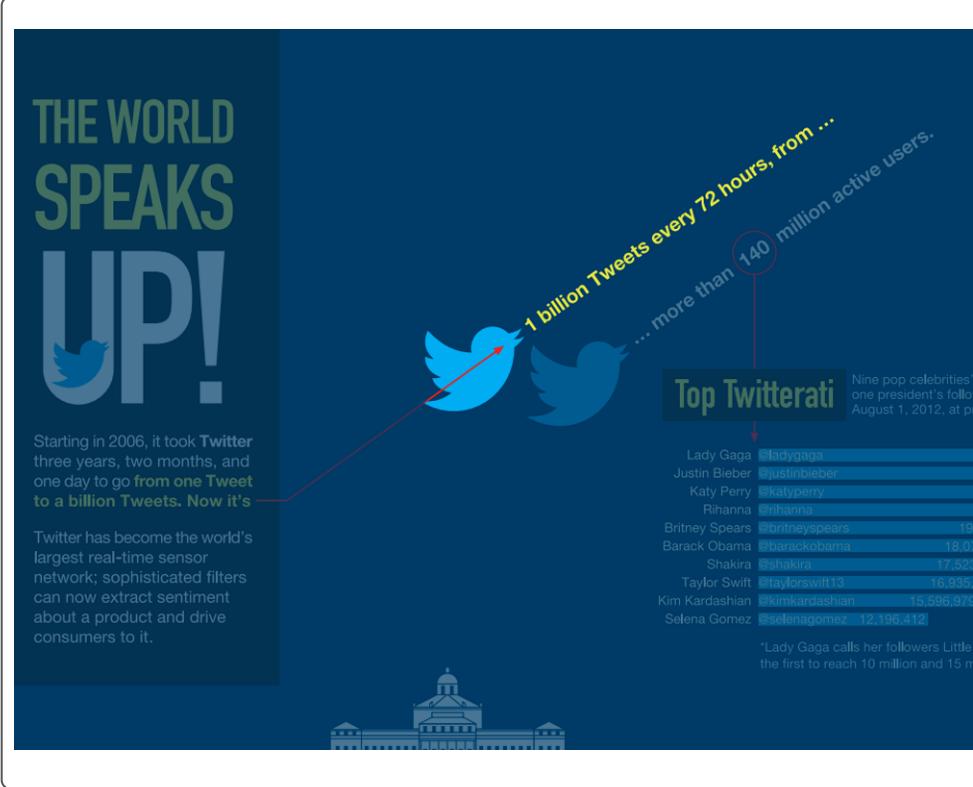
Much of the research centers on how people communicate, why they communicate, and what drives their emotional responses to events in the world and in their own lives. Although the researchers are using Big Data, their efforts attempt to get away from pure numbers for the sake of numbers, and get at much deeper trends in human behavior, stemming from the fact that we are becoming ever more connected.

In 2011, using the entire network of Facebook users, the Data Science Team replicated the famous 1967 "six degrees of separation" experiment.



The Facebook team found that today there are only
4 degrees of separation between us.
(Put it this way:
a friend of your friend
probably knows a friend
of their friend.)





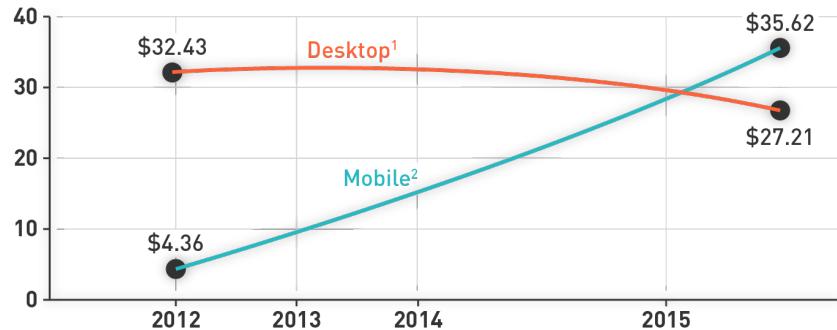
bluefin
LABS



Advertising

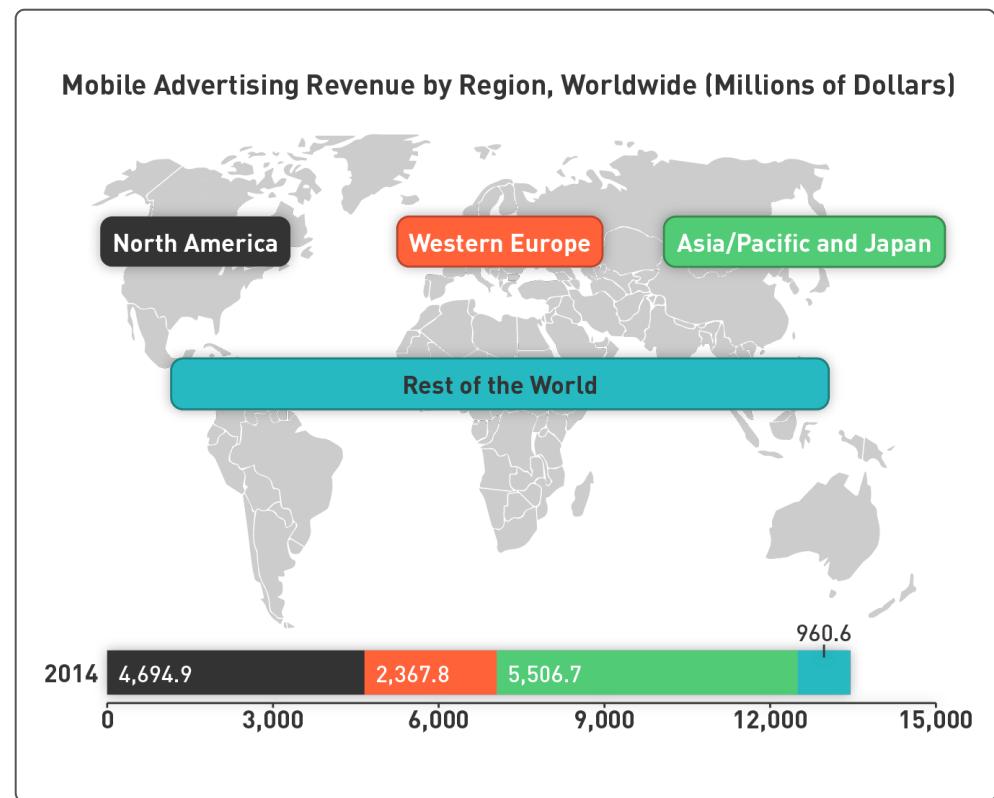
- Represents 2% of GDP in United States since 1900s
- Increased spending online as opposed to on traditional channels

U.S. Digital Ad Spending by Channel, in Billions of Dollars



1- Includes spending primarily on desktop-based ads.

2- Includes classified, display (banners and other, rich media, and video), e-mail, lead generation, messaging-based, and search advertising; ad spending on tablets included.

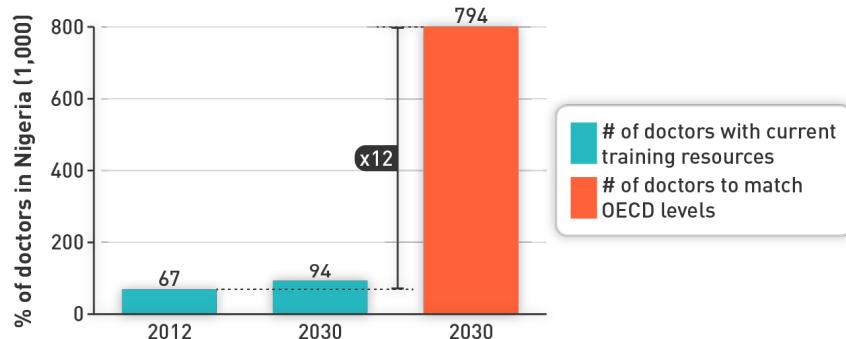


Health Care

Uneven across the world

Imitating Traditional Development Paths Is Impossible for Emerging Economies.

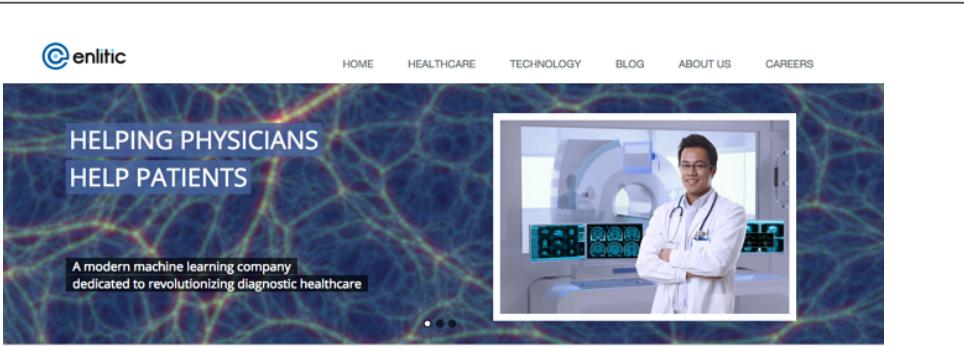
Nigeria would need over 700,000 additional doctors to reach OECD levels by 2030.



Sources: World Bank, WHO, Africa Health Workforce Observatory, BMI, IFC, BCG

- Nigeria: Population 170 million, 67,000 doctors (1:2600 ratio)
- United States: Population 300 million, 1 million doctors (1:300 ratio)

Technology can help address discrepancies in health care.



The screenshot shows the homepage of Enlitic, a machine learning company. The header features the Enlitic logo and navigation links for HOME, HEALTHCARE, TECHNOLOGY, BLOG, ABOUT US, and CAREERS. A main banner with a blue and yellow abstract background displays the text "HELPING PHYSICIANS HELP PATIENTS" and "A modern machine learning company dedicated to revolutionizing diagnostic healthcare". Below the banner is a photo of a doctor in a white coat standing in front of a medical imaging machine. A welcome message "Welcome To Enlitic" is followed by a portrait of Founder and CEO, Jeremy Howard, and a quote from him.

Welcome To Enlitic

 Enlitic uses recent advances in machine learning to make medical diagnostics faster, more accurate, and more accessible. The company's mission is to provide the tools that allow physicians to fully utilize the vast stores of medical data collected today, regardless of what form they are in - such as medical images, doctors' notes, and structured lab tests. To realize this vision, we are building on state-of-the-art deep learning algorithms and partnering with top research hospitals and medical device manufacturers.

"Medical diagnostics is, at its heart, a data problem - turning images, lab tests, patient histories, and so forth into a diagnosis and proposed intervention. Recent applied machine learning breakthroughs, especially using deep learning, have shown that computers can rapidly turn large amounts of data of this kind into deep insights, and find subtle patterns. This is the biggest opportunity for positive impact using data that I've seen in my 20+ years in the field."

— Founder and CEO, Jeremy Howard

Jeremy Howard: Applying machine learning to medical diagnostics

Genomics

- Cost of gene sequencing: \$100,000 (2005) to \$4,000 today
- Near future, \$40 at your local pharmacy

In the United States, 18% of GDP is spent on health care.

**The
Economist****Wireless health care**

When your carpet calls your doctor

The coming convergence of wireless communications, social networking and medicine will transform health care

Apr 8th 2010 | NEW YORK | From the print edition

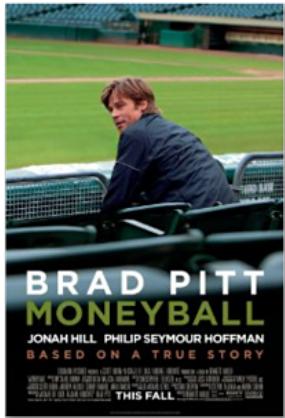


Project of Intel-GE Care Innovations Lab

Home carpet that senses movement, helpful for seniors "aging in place"

Sports

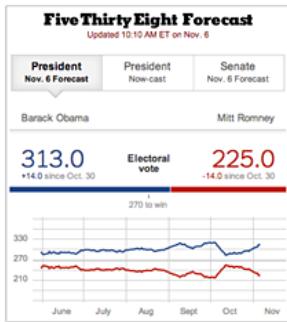
- Oakland Athletics baseball team in Oakland, CA



- Data science used to select best team on small budget

Government

- Barack Obama's victories in 2008 and 2012



- Attributable to effective use of data science

This is the most data driven President we've ever had:

- Created the first set of dashboards at the Federal level to monitor progress on major IT technology investments.
- Established data.gov which hosts over 135,000 data sets (and growing) from the U.S. Government.
- Executive order to ensure that open and machine-readable data is the new default for the government.
- Investing in research and data science to revolutionize how we improve health and treat disease
- Driving privacy for the consumer and ensuring competitiveness through the Big Data report.
- Establishing data driven culture throughout the government with key data personnel at agencies like NIH, Dept of Energy, Commerce, Treasury, Dept. of Transportation, ...

From DJ Patil's Presentation at Strata+Hadoop World

AMP Lab (UC Berkeley)



Office of Science and Technology Policy
Executive Office of the President
New Executive Office Building
Washington, DC 20502

FOR IMMEDIATE RELEASE
March 29, 2012

Contact: Rick Weiss 202 456-6037 rweiss@ostp.eop.gov
Lisa-Joy Zgorski 703 292-8311 lisajoy@nsf.gov

OBAMA ADMINISTRATION UNVEILS "BIG DATA" INITIATIVE: ANNOUNCES \$200 MILLION IN NEW R&D INVESTMENTS

National Science Foundation: In addition to funding the Big Data solicitation, and keeping with its focus on basic research, NSF is implementing a comprehensive, long-term strategy that includes new methods to derive knowledge from data; infrastructure to manage, curate, and serve data to communities; and new approaches to education and workforce development. Specifically, NSF is:

- Encouraging research universities to develop interdisciplinary graduate programs to prepare the next generation of data scientists and engineers;
- Funding a \$10 million Expeditions in Computing project based at the University of California, Berkeley, that will integrate three powerful approaches for turning data into information - machine learning, cloud computing, and crowd sourcing;



Make Money
by working on HITs

HITs - Human Intelligence Tasks - are individual tasks that you work on. [Find HITs now.](#)

As a Mechanical Turk Worker you:

- Can work from home
- Choose your own work hours
- Get paid for doing good work

Find an interesting task → Work → Earn money

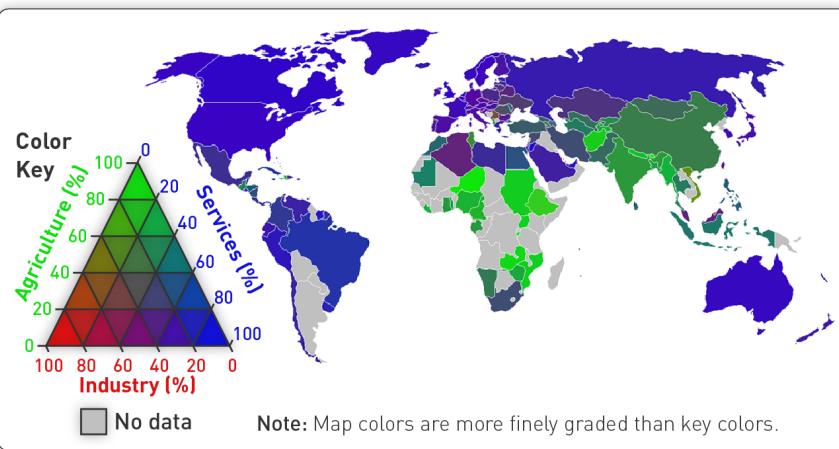
or [learn more about being a Worker](#)

[Find HITs Now](#)

The diagram illustrates the workflow: finding an interesting task leads to work, which results in earning money. It also includes a sidebar text about HITs and a link to learn more about being a worker.

Micro work-based framework connecting short tasks with workers

Labor Force by Occupation



First world: Service oriented, can leverage machine learning

Summary

- Many big data problems where humans play a central role
- Next: Contrast with machine-centric big data

We are at a tipping point—humans are no longer at the center of the data universe.

The Internet of Things (IoT)



More data manufactured by machines—servers, cell phones, GPS-enabled cars—than by people

Advent of IoT Era

- Machine-to-machine learning (M2M)
 - Attach sensors to things animate and inanimate, and network these
 - Stream data
 - By 2020, over 26 billion things in network (majority M2M)
- Has evolved from convergence of wireless technologies, micro-electromechanical systems (MEMS), and the Internet
- Is a scenario in which objects, animals, or people are provided with:
 - Unique identifiers
 - The ability to transfer data over a network without human-to-human or human-to-computer interaction

Japan Earthquake



- Fourteen seconds before, all Shinkansen trains (bullet trains) and plants shut down
- Mesh of sensors from connected computers to form an early warning system

Tracking Nature in the Wild



- Sea lions instrumented with GPS tracking system
- Underwater hub picks up signals

Serengeti



- Animals tracked through attached GPS systems or through sensors (motion detectors and cameras)
- Work of Tanya Berger-Wolf, 2014

Challenges

- Battery life, storage, network connectivity
- Privacy and security

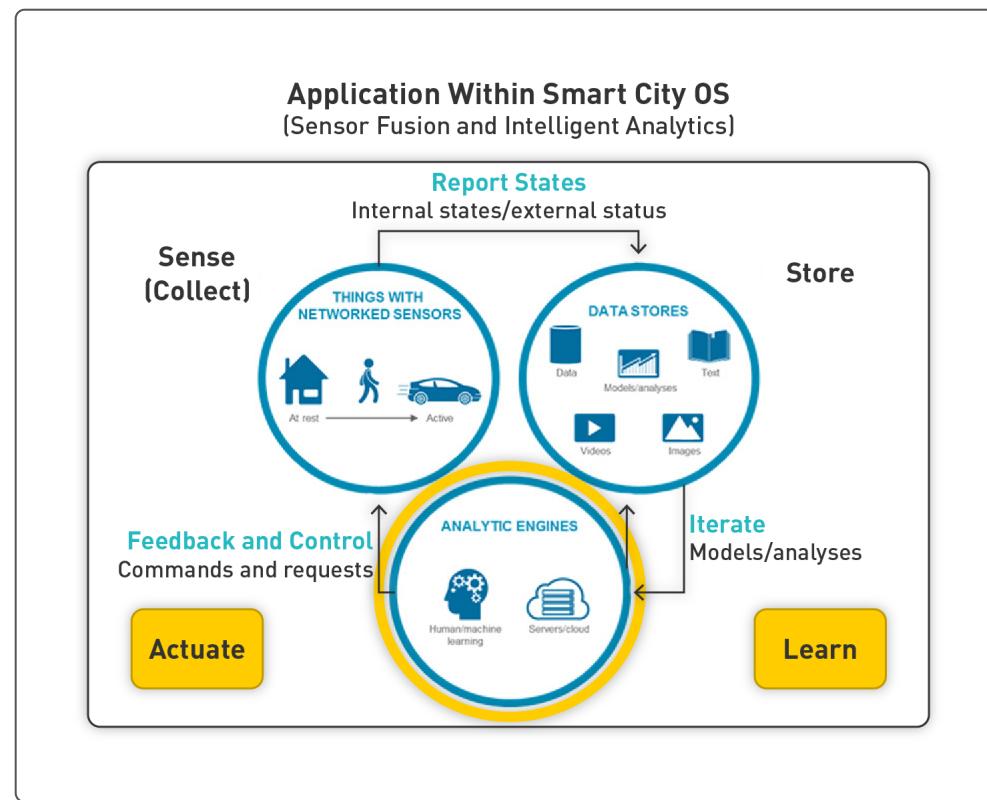
Smart Homes



- By 2020, 1 billion smart meters
- One megabyte per device daily (from polling 1,000 times) → 1 petabyte of data daily

Smart Cities

- Involves capturing telemetry data from IoT and citizens (cell phones)
- Affects all sectors of society, including:
 - Government services
 - Transport and traffic management
 - Energy
 - Health care
 - Water
 - Waste



Autonomous Vehicles

- Driverless cars (Google)
- Drones

Big data is not just about Web 2.0/Social, but IoT also.

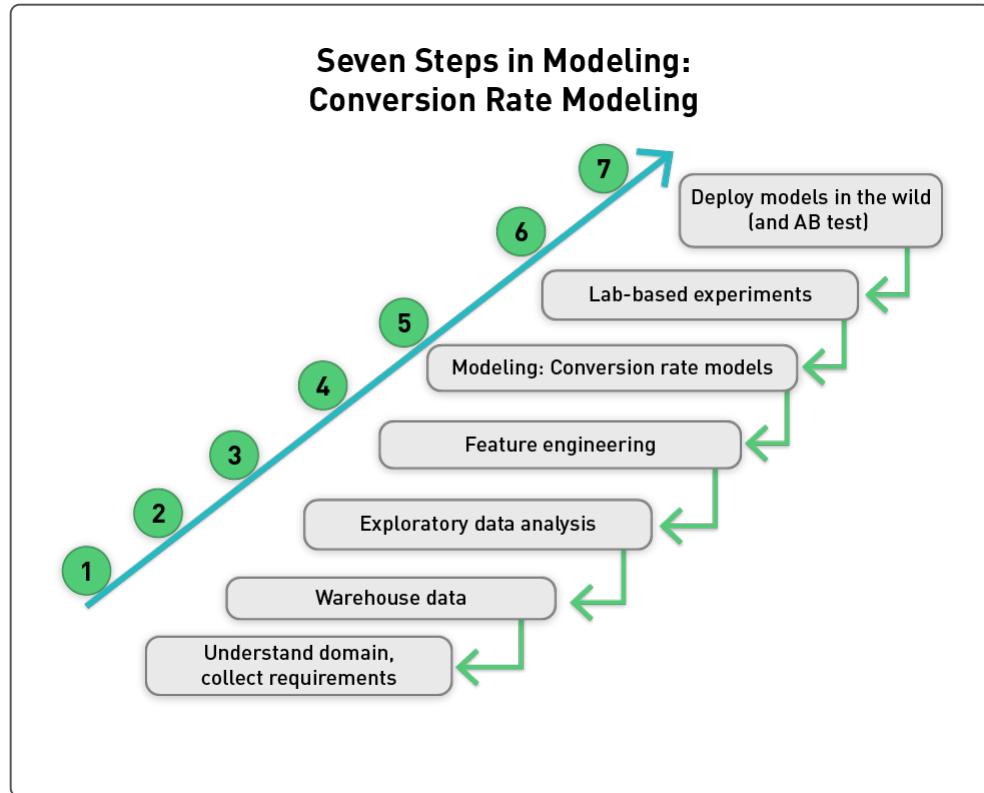


THE cloudera SESSIONS
your inside track to big data

By 2020, over 26 billion devices will make up the IoT.

Introduction

- New opportunities and business needs generated by big data applications
- Massive opportunity and challenge for machine learning
- Machine learning does not exist in a vacuum
 - Requires a sophisticated ecosystem

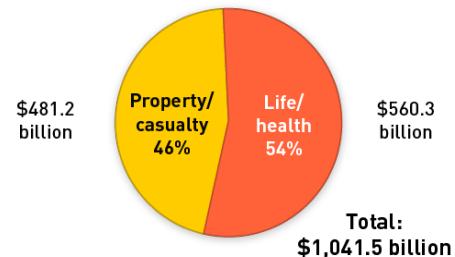


Where does time go in large-scale machine learning?

Eighty percent of the work in any data project is preparing the data.

Example: U.S. Insurance Industry

U.S. Property/Casualty and Life/Health Insurance Premiums, 2013



- Net premiums totaled \$1 trillion in 2013
 - Almost half property and casualty (auto, home, commercial)
- Highly competitive space, leading to development of new programs
 - E.g., special rates for good drivers
 - Score computed through accident history over, say, five years
 - Static, annual, small-scale, batch computation

Progressive's Snapshot Product

- "Bad driver" surveillance using telematics
- Via small box plugged into steering wheel that records, sends information
- Features:
 - How many miles driven
 - How many miles driven between midnight and 4 a.m.
 - Use of sharp or gentle braking
- Data used to analyze driving patterns, risk to insurance company
- Potential reduced premiums for user

Privacy

- Huge data exhaust being generated from digital and online activities
- Data not transitory but permanent
- Caution needed with storage, accessibility

Challenges

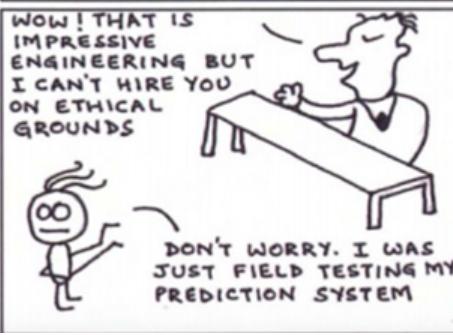
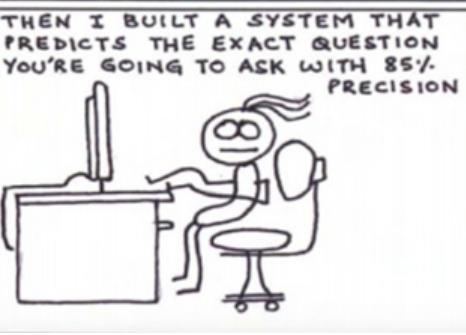
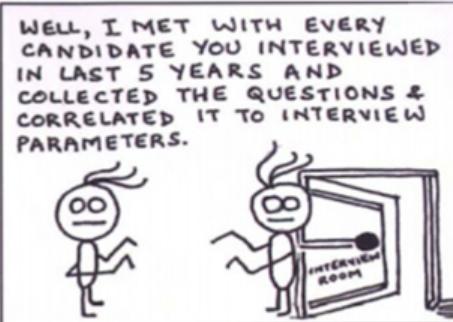
- International security
- Privacy

Pacemaker: Privacy vs. Ownership

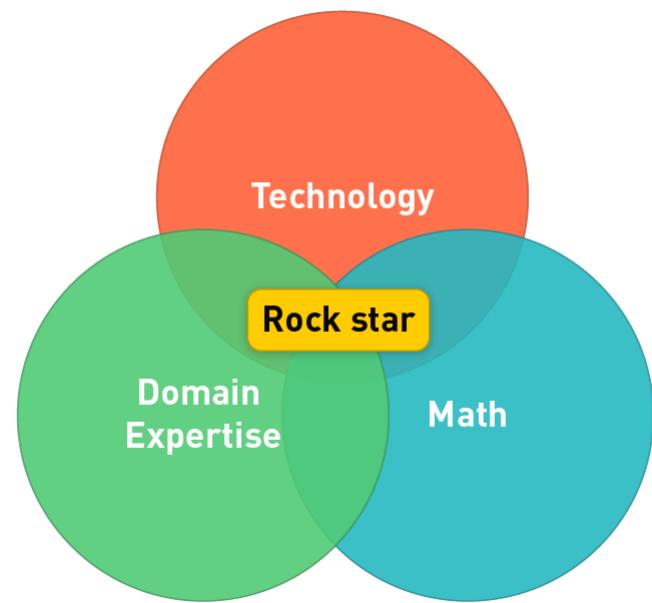


Who owns the data we generate?

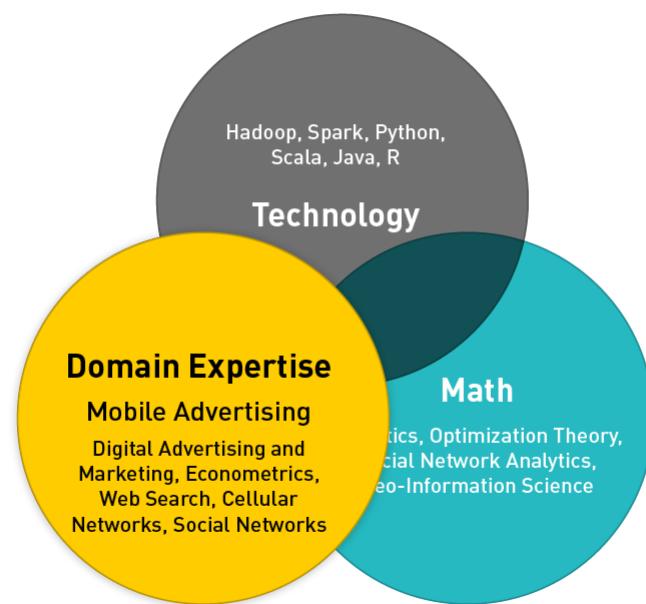
When you interview a data scientist...

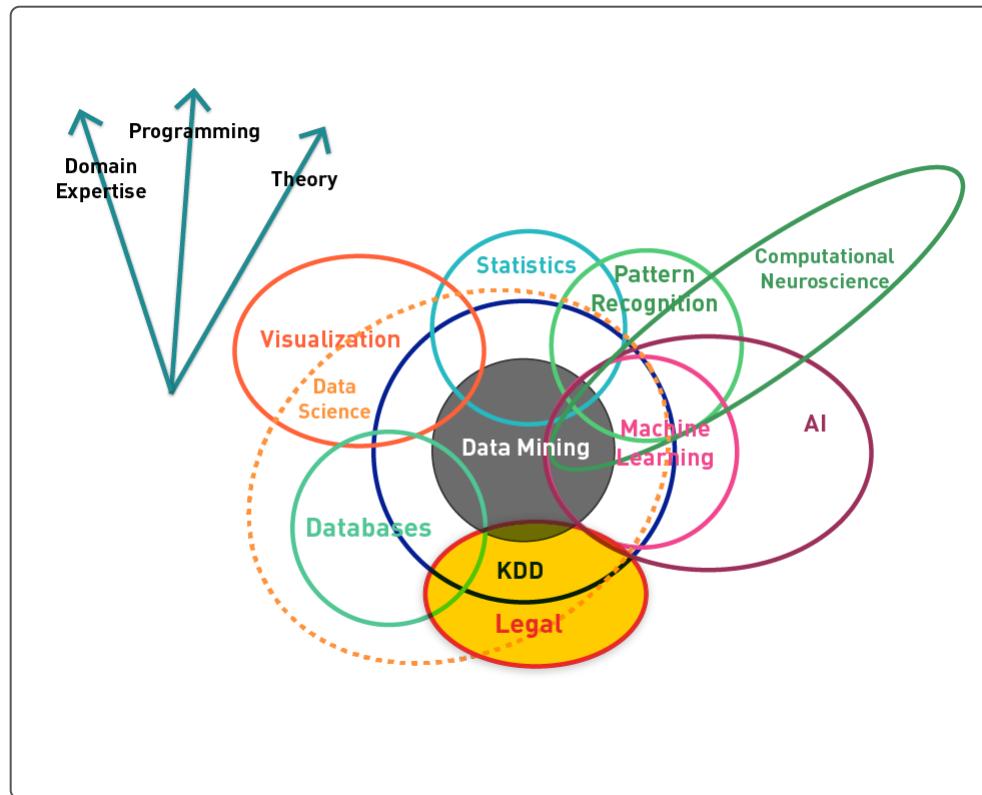


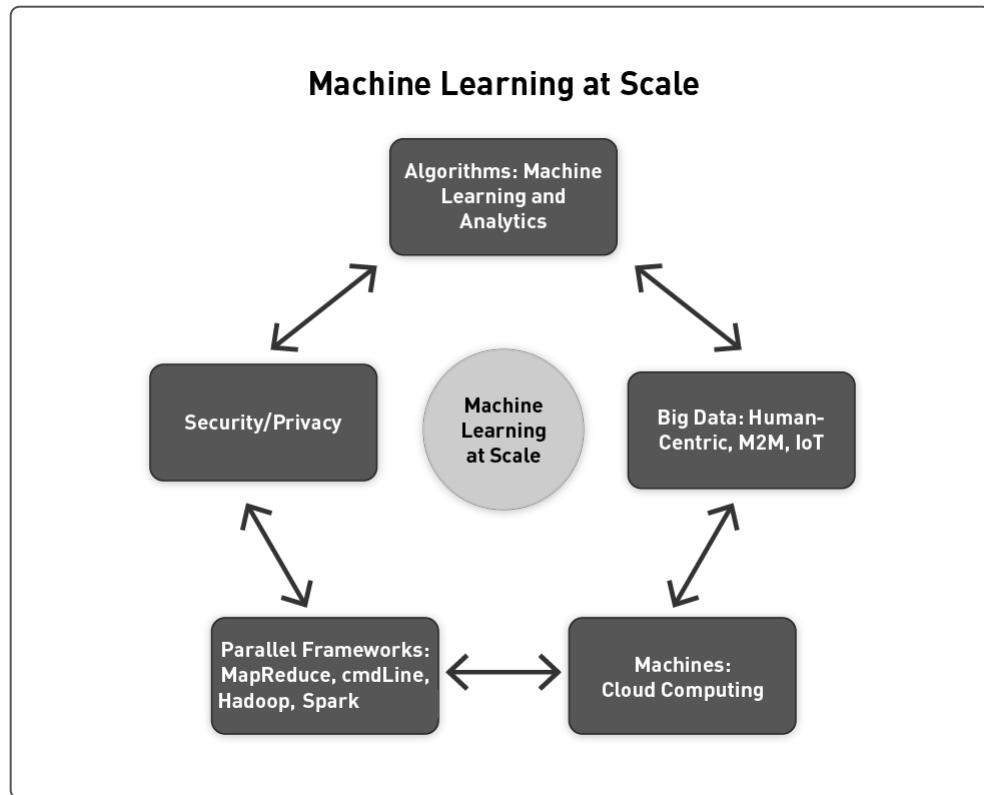
Three Skills Sets for Large Scale Machine Learning



Three Skills Sets in Detail









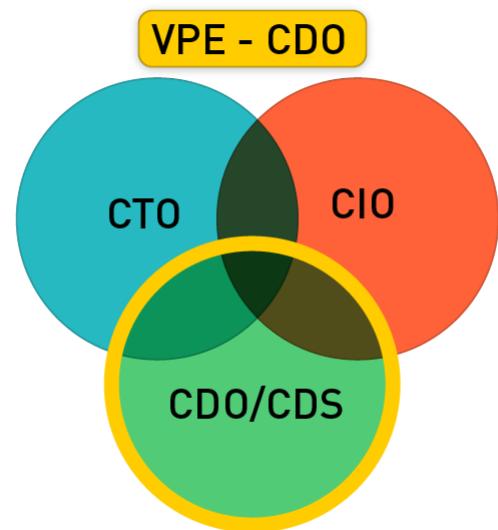
“Swiss army knife of the 21st century”

—MediaGuardian Innovation Awards

Data Science (DS) Team

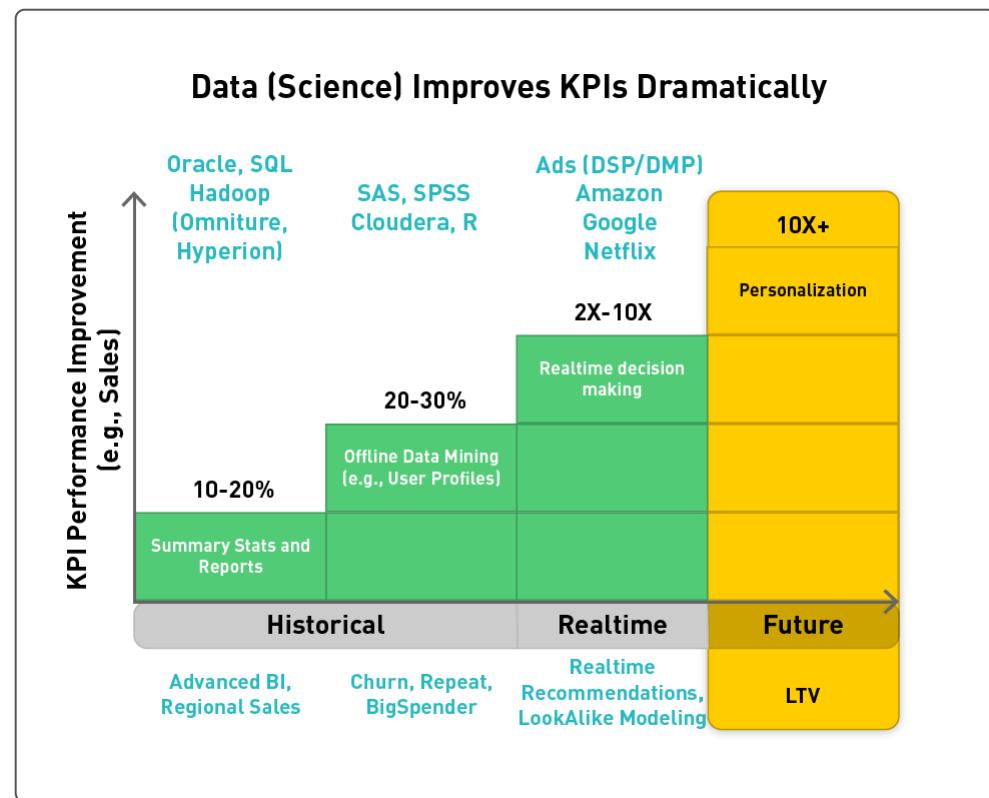
- Experts in:
 - Technology
 - Math
 - Domain
 - Legal (privacy, security)
- Project managers
- Research scientists
- Communications experts

Evolution of Data Science Functional Role



The challenge is to leverage data at scale ...

... as we move from era of business intelligence to data science.



Data Scientists in Huge Demand



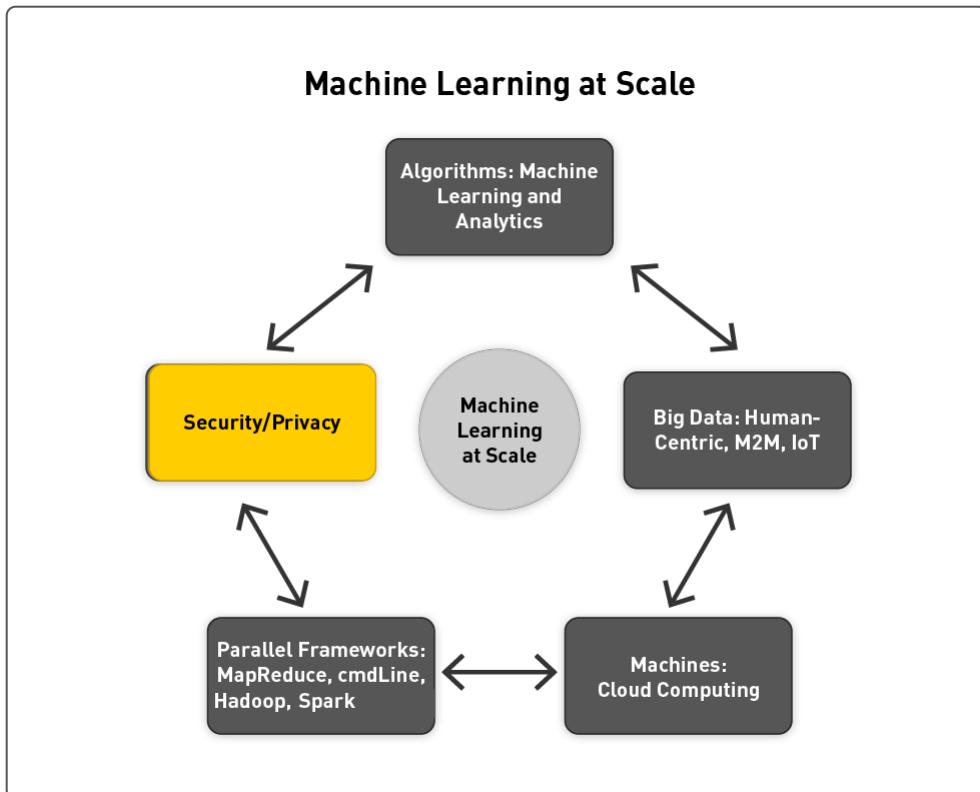
There will be over 2 million open positions in the United States.

After This Class ...

- Role
 - Individual contributors: R&D, r&D, R&d, D
 - Managers and leaders
- Focus
 - Research
 - Continue studies, get PhD
 - Teach
 - Conduct theoretical and applied research

After This Class ... (cont.)

- Focus
 - Infrastructure development: Build
 - Development
 - Architects of big data pipelines and large-scale machine learning
 - Full-stack people
 - Builders of apps (on a fully supported framework)



Large-Scale Machine Learning

- Needs supporting infrastructure and ecosystem
- Parallel computing frameworks needed

MapReduce Framework

- Allows us to divide and conquer huge data problems
- Four frameworks
 - Command-line-based approach
 - Hadoop
 - ~~MRJob~~ (This is no longer covered in the course)
 - Spark

"Embarrassingly Parallel" Problems

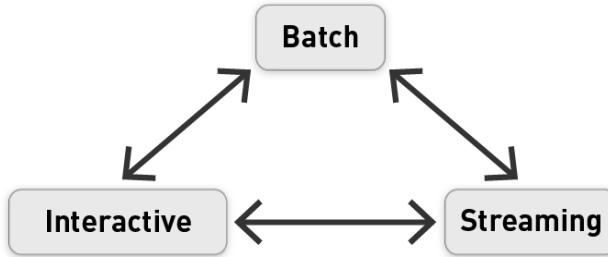
- Most machine learning problems
- Divided into subproblems, which:
 - Require little or no communication
 - Are easily distributed
 - Have linear computational flow and scale

Most machine learning algorithms are amenable to parallel computation.

Framework Requirements

- Scalable
- Fault tolerant
- Iterative
- Interactive

Few frameworks provide all these aspects



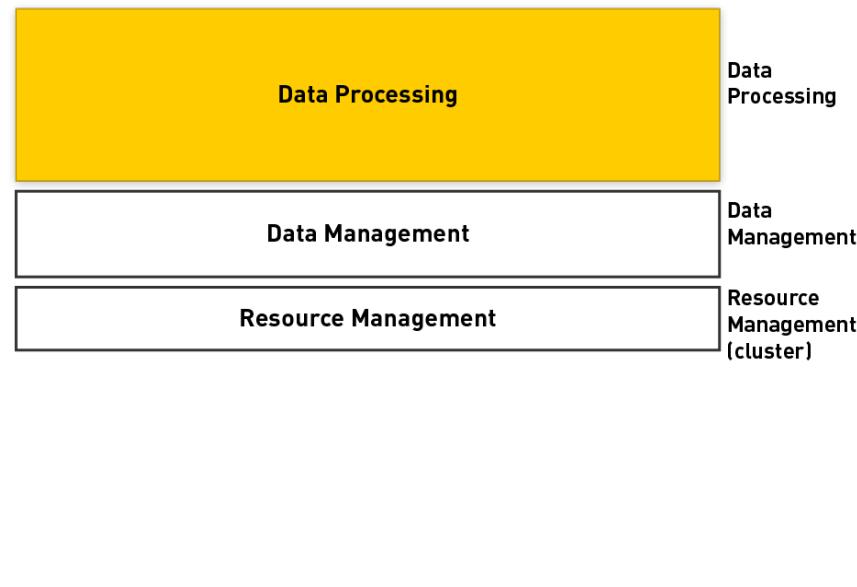
- Also compatible with existing open-source ecosystem

Four MapReduce Frameworks

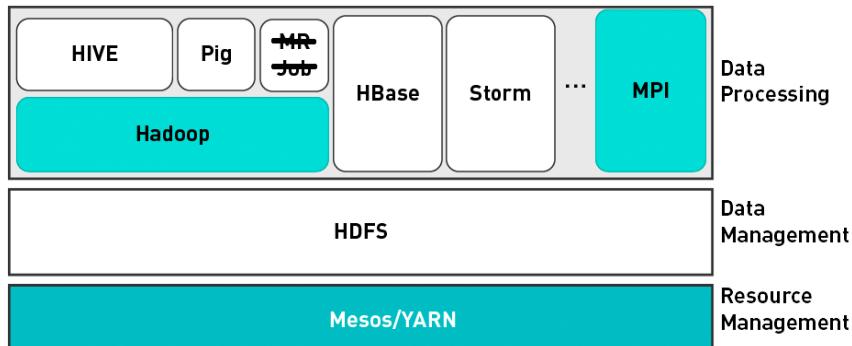
In increasing order of sophistication and ease of use

- Command-line-based approach
- Hadoop: Storing and manipulating data
- ~~MRJob~~ (This is no longer covered in the course)
- Spark:
 - Memory backed
 - Integrated framework for variety of activities
 - Specialized libraries for machine learning, graph processing, data analytics

Data Analytics Stack

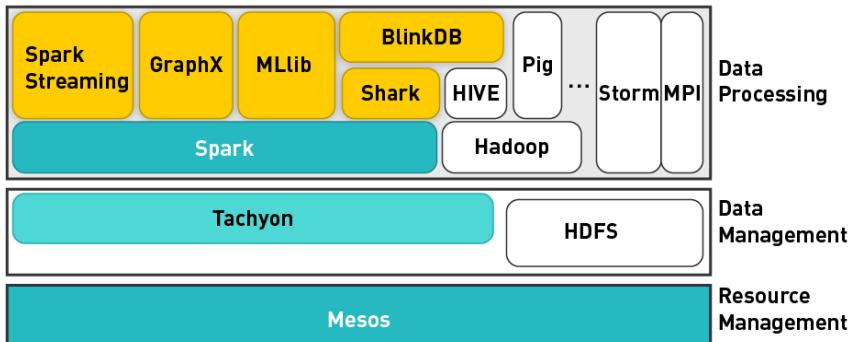


Populating the Stack

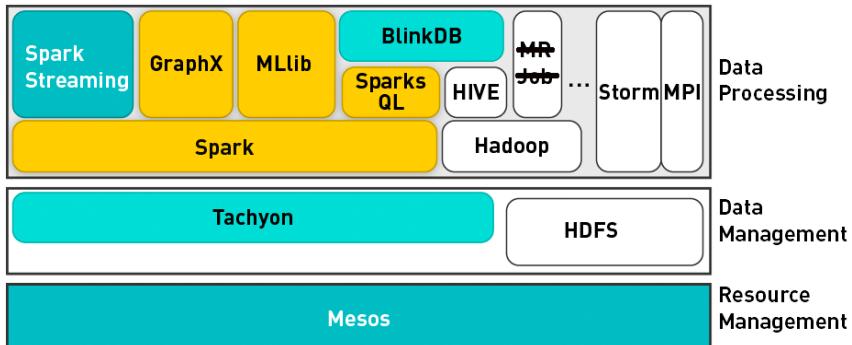


Mesos/YARN: Allow multiple frameworks on single cluster

Populating the Stack: Spark SQL



Class Phases



- Phase 0: Command line
- Phase 1: Hadoop/HDFS
- Phase 2: MRJob (This is no longer covered in the course)
- Phase 3: Spark

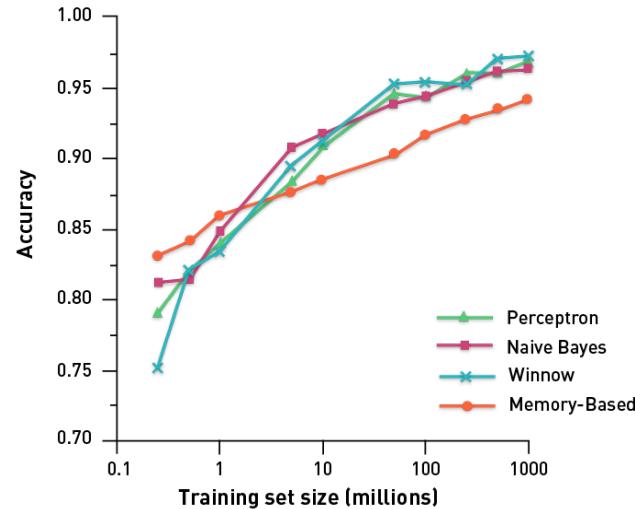
Machine Learning Algorithms

- **Supervised machine learning** (convex optimization, gradient descent, linear regression, decision trees, ensembles of models, support vector machines)
- **Unsupervised** (expectation maximization, matrix multiplication, alternating least squares)
- **Graphs** (random walks, PageRank, graph search algorithms such as breadth-first search, shortest path)
- **Hybrid algorithms** (supervised machine learning and random walks)
- **Applications** (digital advertising, social media, health care, e-commerce, entertainment, metrics, statistics)

The question: More data scientists or more data?

Machine Learning and Data Study

- Filling in confusable words in sentences
 - E.g., *For breakfast I ate ___ eggs. Use {to, two, too}.*



- Conclusion: More data leads to 10–20% boost in performance

Supervised Classification in a Nutshell

Given $D = \{(x_i, y_i)\}_i^n$

A diagram showing a sparse feature vector x_i and a label y_i . A red bracket under x_i is labeled '(sparse) feature vector'. A red arrow points from y_i to its position in the set, labeled 'label'.

Induce $f : X \rightarrow Y$ s.t. loss is minimized

$$\text{empirical loss} = \frac{1}{n} \sum_{i=0}^n \ell(f(x_i), y_i)$$

A diagram showing the empirical loss formula. A red arrow points from the term $\ell(f(x_i), y_i)$ to its position in the sum, labeled 'loss function'.

Consider functions of a parametric form:

$$\arg \min_{\theta} \frac{1}{n} \sum_{i=0}^n \ell(f(x_i; \theta), y_i)$$

A diagram showing the parametric form of the loss function. A red arrow points from the term $f(x_i; \theta)$ to its position in the function, labeled 'model parameters'.

Key insight: machine learning as an optimization problem!
(closed form solutions generally not possible)



Gradient Descent

$$w^{(t+1)} = w^{(t)} + \gamma^{(t)} \frac{1}{n} \sum_{i=0}^n \nabla l(f(x_i; \theta^{(t)}), y_i)$$

“batch” learning: update model after considering all training instances

Stochastic Gradient Descent (SGD)

$$w^{(t+1)} = w^{(t)} + \gamma^{(t)} \nabla l(f(x; \theta^{(t)}), y)$$

“online” learning: update model after considering each (randomly-selected) training instance

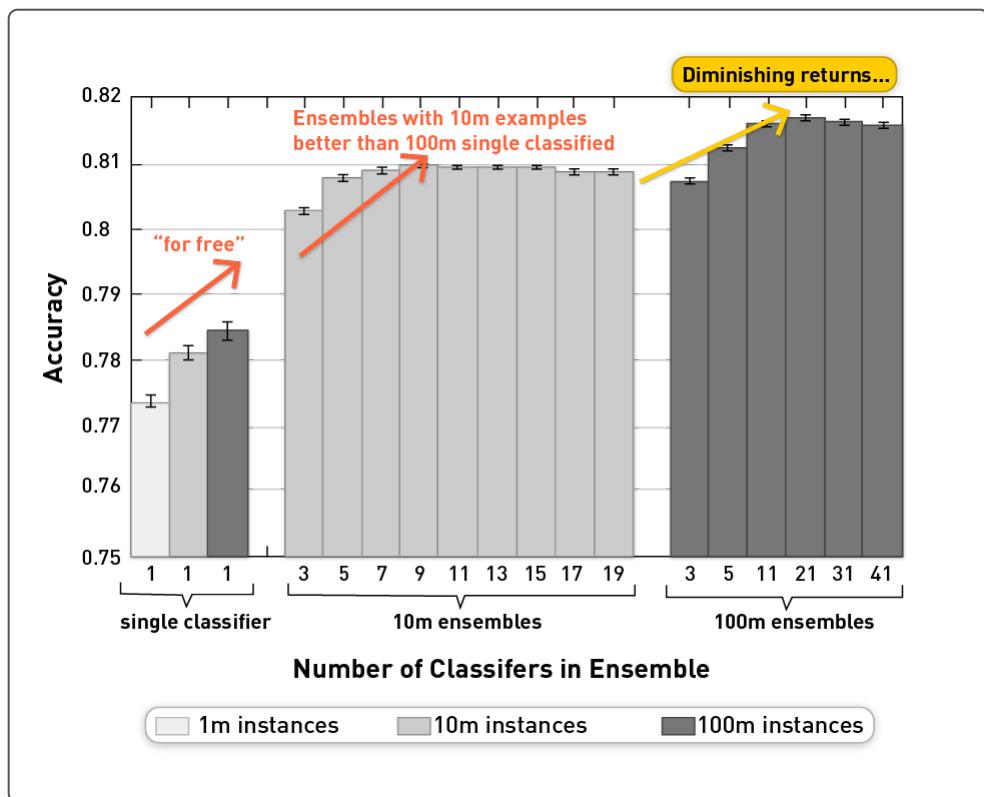
In practice... just as good!

Solves the iteration problem!

What about the single reducer problem?



**Ensembles of single
models can help augment our performance.**



More Data or More Data Scientists?

- Trade-off between bias and variance
- Bias from data scientists
- Variance from more data
 - With more data can reduce variance

Bias-Variance

- Empirical studies: More data leads to big improvements
- Gigs of training data to TBs or PBs: Improved performance
- Formal perspective: Bias-variance trade-off

Machine Learning Objectives

- Minimize error
- Minimize loss function
- Reduce model complexity
- Better generalization
- E.g., mean squared error for regression

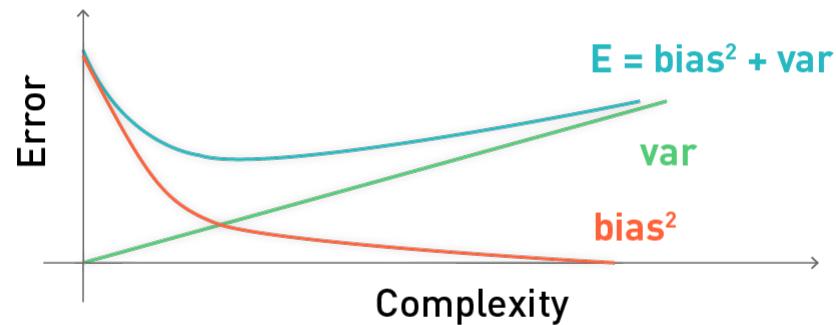
Loss: Irreducible and Reducible Error

- Irreducible:
 - Inherent uncertainty
 - Associated with natural variability in system
 - E.g., noisy sensors, incorrectly documented data
- Reducible:
 - Can and should be minimized to maximize accuracy

Reducible Error

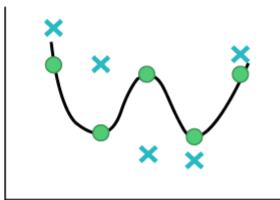
- Includes both "error due to squared bias" and "error due to variance"
- Goal: Simultaneously reduce bias and variance as much as possible in order to obtain as accurate a model as is feasible.
- Trade-offs in selecting models:
 - Flexibility and complexity
 - Selecting appropriate training sets

Complexity of the Model

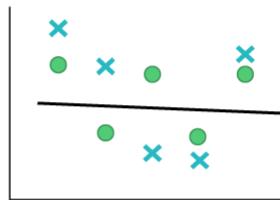


Usually, bias is a decreasing function of complexity, while variance is an increasing function of complexity.

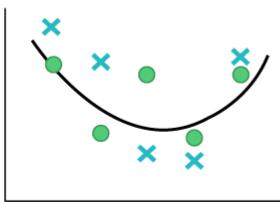
Bias-Variance Trade-Off in Model Selection in Simple Problem



(a) High variance/low bias.
Fourth-order polynomial ($p = 5$).



(b) Low variance/high bias.
First-order polynomial ($p = 2$).



(c) Balanced variance and bias.
Minimum MSE.
Second-order polynomial ($p = 3$).

● Data points for fitting
✖ Typical new data points

Error Due to Squared Bias

- Definition: Amount by which expected model prediction differs from the true value, over the training data
- Bias is introduced at model selection
- Repeat model-building process (through resampling) to obtain average of prediction values
- If average prediction values are substantially different than true value, bias will be high

Error Due to Variance

- Definition: Amount by which prediction over one training set differs from expected predicted value over all training sets
- Repeat modeling, measure inconsistency over different training data sets
- Variance measures how inconsistent predictions are from one another, over different training sets, not whether or not they are accurate (Manning et al., 2008)

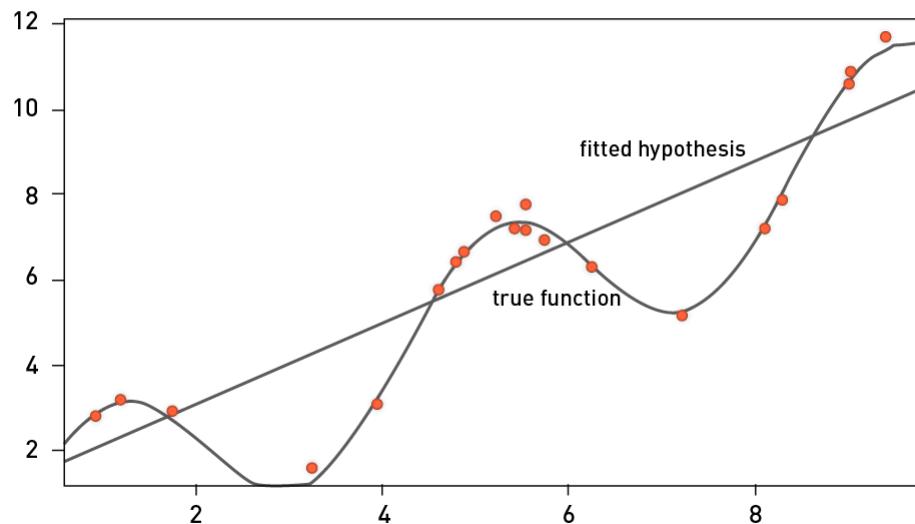
Bias-Variance Analysis in Regression

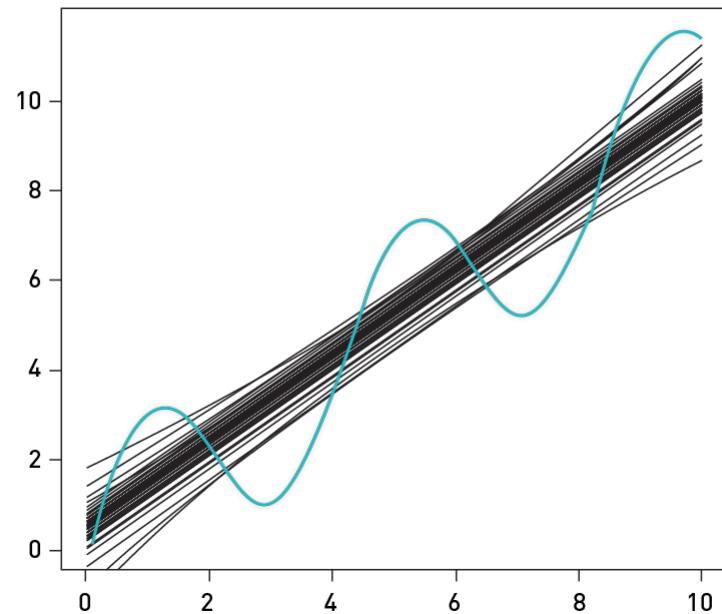
- True function: $y = f(x) + \varepsilon$
 - Where ε is normally distributed with zero mean and standard deviation σ .
- Given a set of training examples, $\{(x_i, y_i)\}$, we fit a hypothesis $h(x) = wx + b$ to the data to minimize the squared error.

$$\sum_i [y_i - h(x_i)]^2$$

Fit a Linear Hypothesis: 20 Points

$$y = x + 2 \sin(1.5x) + N(0, 0.2)$$



50 Fits (20 Examples Each)

Variance, Bias, and Noise

- Variance:

$$\mathbb{E}h(x^*) - \mathbb{E}[h(x^*)]^2$$

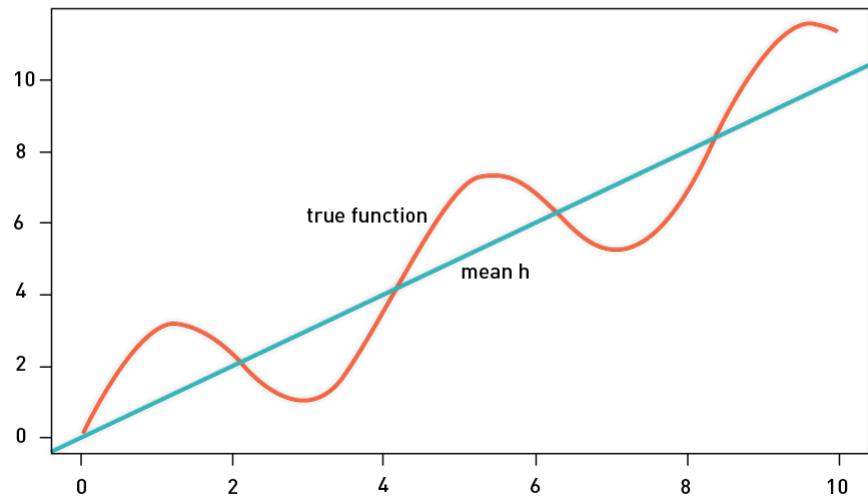
- Bias:

$$\mathbb{E}[h(x^*)] - f(x^*)$$

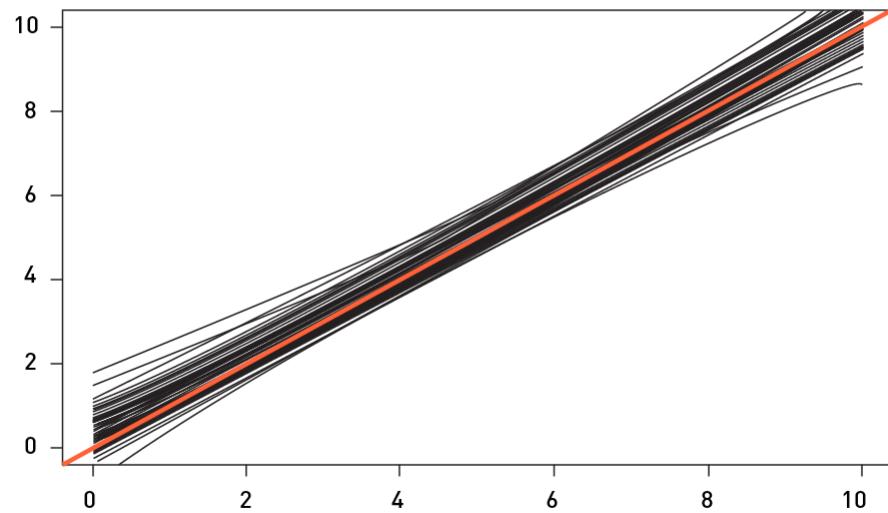
- Noise:

$$\mathbb{E}[(y^* - f(x^*))^2] = \mathbb{E}[\varepsilon^2] = \sigma^2$$

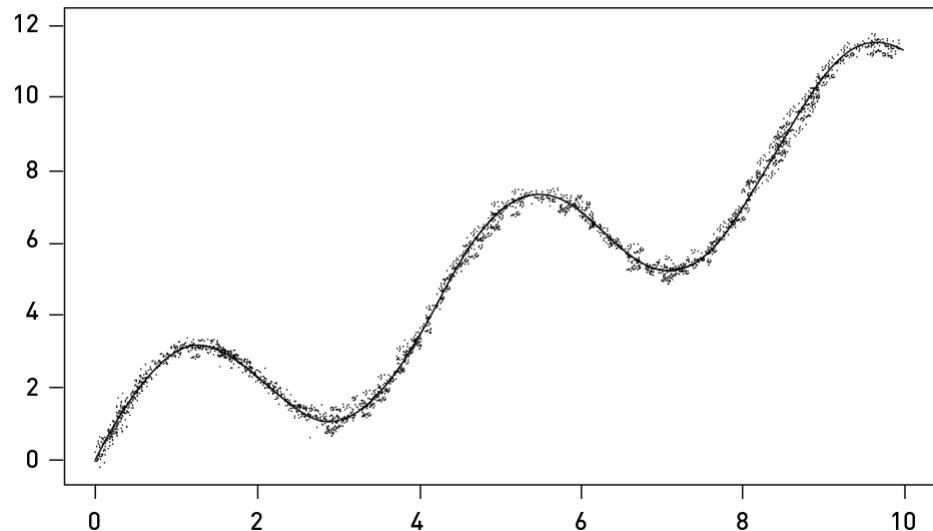
Bias

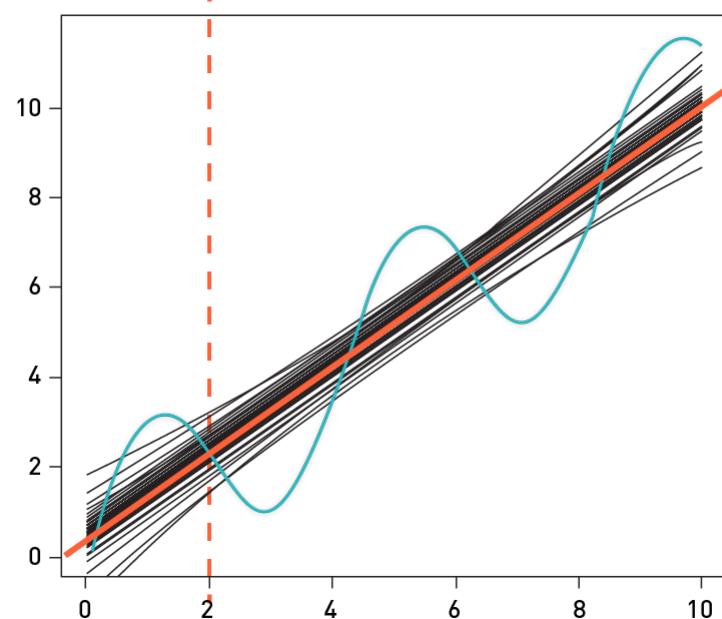


Variance

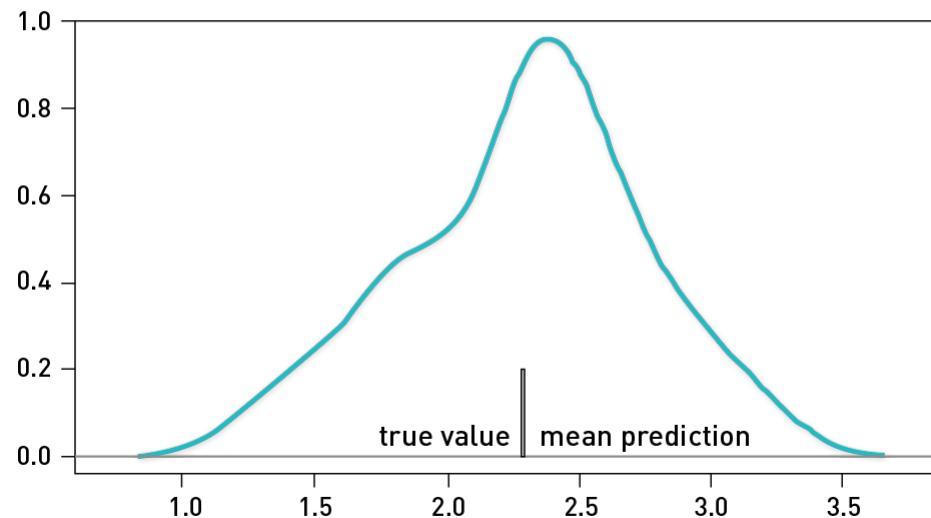


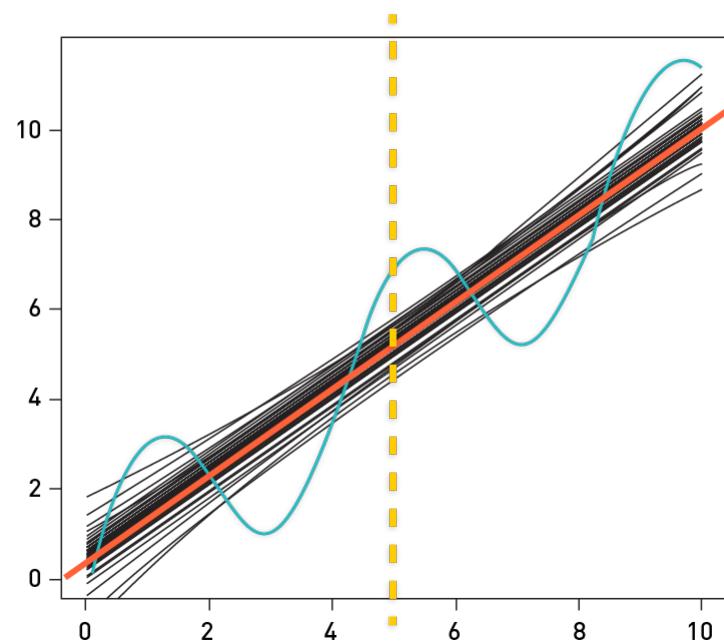
Noise



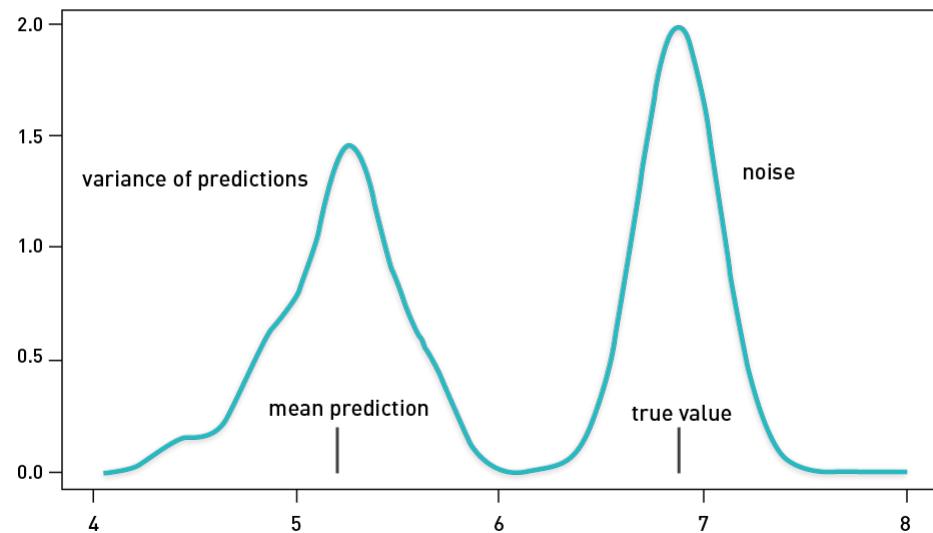
Distribution of Predictions at $x = 2.0$ 

Distribution of Predictions at $x = 2.0$



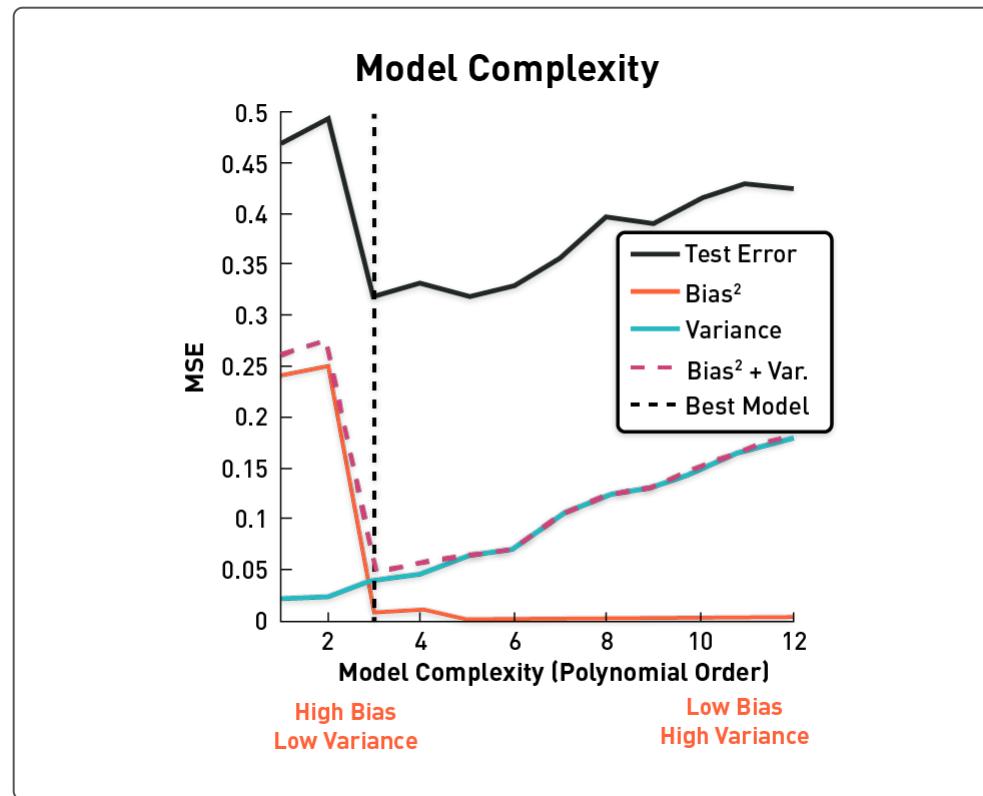
Distribution of Predictions at $x = 5.0$ 

Distribution of Predictions at $x = 5.0$



Measuring Bias and Variance

- Simulate multiple training sets by bootstrap replicates.
- Construct B bootstrap replicates of S (e.g., $B = 200$).
 - S_1, \dots, S_B
- Apply learning algorithm to each replicate S_b to obtain hypothesis h_b .
- Let $T_b = S \setminus S_b$ be the data points that do not appear in S_b .
- Compute predicted value $h_b(x)$ for each x in T_b .



Trade-off between bias and variance: Best to minimize both at the same time

Bias-Variance Decomposition

- Can be extended to classification problems
- Pedro Domingos (2000a; 2000b): Developed unified decomposition that covers both regression and classification

Bias-Variance Trade-Off Code: Articles

Excellent articles with code in matlab for bias-variance analysis of squared error loss:

- [Model Selection: Underfitting, Overfitting, and the Bias-Variance Tradeoff](#)
- [Ask a Data Scientist: The Bias vs. Variance Tradeoff](#)
- [polyfit: Polynomial curve fitting](#)

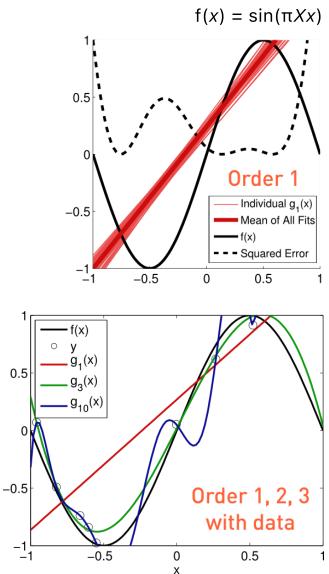
Bias-Variance Order: 1, 2, 3 Polynomials

```

xGrid = linspace(-1, 1, 100);
1 % FIT MODELS TO K INDEPENDENT DATASETS
2 K = 50;
3 for iS = 1:K
4   ySim = f(x) + noiseSTD*randn(size(x));
5   for jD = 1:numel(degree)
6     % FIT THE MODEL USING polyfit.m
7     thetaTmp = polyfit(x,ySim,degree(jD));
8     % EVALUATE THE MODEL FIT USING polyval.m
9     simFit{jD}(iS,:) =
polyval(thetaTmp,xGrid);
10    end
11  end
12
13 % DISPLAY ALL THE MODEL FITS
14 h = [];
15 for iD = 1:numel(degree) For polynomial = iD
16   figure(iD+1)
17   hold on
18   % PLOT THE FUNCTION FIT TO EACH DATASET
19   for iS = 1:K
20     h(1) = plot(xGrid,simFit{iD}
(iS,:),'color',
brighten(cols(iD,:),.6));
21   end
22   % PLOT THE AVERAGE FUNCTION ACROSS ALL FITS
Average predication of xGrid
23   h(2) = plot(xGrid,mean(simFit{iD}),
'color',cols(iD,:),'Linewidth',5);
24   % PLOT THE UNDERLYING FUNCTION f(x)
25   h(3) = plot(xGrid,f(xGrid),'color','k','Linewidth',3);
26   % CALCULATE THE SQUARED ERROR AT EACH POINT, AVERAGED ACROSS ALL DATASETS
BIAS
27   squaredError = (mean(simFit{iD})-f(xGrid)).^2;
28   % PLOT THE SQUARED ERROR
29   h(4) = plot(xGrid,squaredError,'k--','Linewidth',3);
30   uistack(h(2),'top')
31   hold off
32   axis square
33   xlim([-1 1])
34   ylim([-1 1])
35   legend(h,{sprintf('Individual g_{%d}(x)',degree(iD)),
'Mean of All Fits','f(x)','Squared Error'},'Location','WestOuts...
36   title(sprintf('Model Order=%d',degree(iD)))
37 end

```

Our original goal was to
approximate $f(x)$, not the data
points, per se.

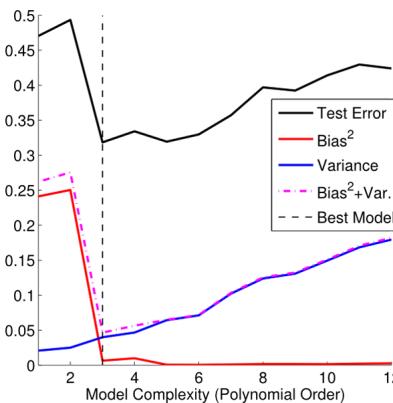


Bias-Variance Order: 1 to 12 Polynomials

```

:
7 % # INITIALIZE SOME VARIABLES
8 xGrid = linspace(-1,1,N);
9 meanPrediction = zeros(K,N);
10 thetaHat = {};
11 x = linspace(-1,1,N);
12 x = x(randperm(N));
13 for iS = 1:K % LOOP OVER DATASETS
14 % CREATE OBSERVED DATA, y
15 y = f(x) + noiseSTD*randn(size(x));
16
17 % CREATE TRAINING SET
18 xTrain = x(1:nTrain);
19 yTrain = y(1:nTrain);
20
21 % CREATE TESTING SET
22 xTest = x(nTrain+1:end);
23 yTest = y(nTrain+1:end);
24
25 % FIT MODELS
26 for jD = 1:nPolyMax
27
28 % MODEL PARAMETER ESTIMATES
29 thetaHat{jD}(iS,:) = polyfit(xTrain,yTrain,jD);
30
31 % PREDICTIONS
32 yHatTrain{jD}(iS,:) = polyval([thetaHat{jD}(iS,:)],xTrain); % TRAINING
SET
33 yHatTest{jD}(iS,:) = polyval([thetaHat{jD}(iS,:)],xTest); % TESTING SET
34
35 % MEAN SQUARED ERROR
36 trainErrors{jD}(iS) = mean((yHatTrain{jD}(iS,:) - yTrain).^2); %
TRAINING
37 testErrors{jD}(iS) = mean((yHatTest{jD}(iS,:) - yTest).^2); % TESTING
38 end
39 end
40 % Test Error=Variancexi+Biasxi+irreducibleError
41
42 % CALCULATE AVERAGE PREDICTION ERROR, BIAS, AND VARIANCE
43 for iD = 1:nPolyMax
44 trainError(iD) = mean(trainErrors{iD});
45 testError(iD) = mean(testErrors{iD});
46 biasSquared(iD) = mean((mean(yHatTest{iD})-f(xTest))x^2); Avg(Bias (xi))
47 variance(iD) = mean(var(yHatTest{iD},1)); Avg(Variance (xi))
48 end
49 [~,bestModel] = min(testError);

```



Model Selection

- Trade-off between bias and variance
- Practical method for selecting model:
 - Minimize error
 - Function of model complexity

Effect of Algorithm Parameters

- K-Nearest Neighbor: Increasing k typically increases bias, reduces variance
- Decision trees of depth D: Increasing D typically increases variance, reduces bias
- Radial basis function (RBF) support vector machines (SVM) with parameter σ : Increasing σ increases bias, reduces variance

Bagging tends to reduce variance without changing bias.

Wrap-Up: Bias-Variance Trade-Off

- Bias and variance of an estimator are related to squared prediction error.
- These concepts can be applied to classification problems.
- An optimal estimator will have both low variance and low bias.

Parallel Sort

```
1  #! /bin/ksh
2
3  MAX_LINES_PER_CHUNK=1000000
4  ORIGINAL_FILE=$1
5  SORTED_FILE=$2
6  CHUNK_FILE_PREFIX=${ORIGINAL_FILE}.split.
7  SORTED_CHUNK_FILES=${CHUNK_FILE_PREFIX}*.sorted
8
9  usage () {
10    echo Parallel sort
11    echo usage: psort file1 file2
12    echo Sorts text file file1 and stores the output in file2
13    echo Note: file1 will be split in chunks up to
14    echo $MAX_LINES_PER_CHUNK lines
15    echo and each chunk will be sorted in parallel
16  }
17
18 # test if we have two arguments on the command line
19 if [ $# != 2 ]
20 then
21   usage
22   exit
23 fi
24
25 #Cleanup any leftover files
26 rm -f ${SORTED_CHUNK_FILES} > /dev/null
27 rm -f ${CHUNK_FILE_PREFIX}* > /dev/null
28 rm -f ${SORTED_FILE}
29
30 #Splitting ${ORIGINAL_FILE} into chunks ...
31 split -l ${MAX_LINES_PER_CHUNK} ${ORIGINAL_FILE}
${CHUNK_FILE_PREFIX}
32
33 for file in ${CHUNK_FILE_PREFIX}*
34 do
35   sort $file > ${file}.sorted &
36 done
37 wait
38
```

```

39 #Merging chunks to $SORTED_FILE ...
40 sort -m $SORTED_CHUNK_FILES > $SORTED_FILE
41
42 #Cleanup any leftover files
43 rm -f $SORTED_CHUNK_FILES > /dev/null
44 rm -f $CHUNK_FILE_PREFIX* > /dev/null

```

General Framework for Mapping and Reducing

```

1  %%writefile cmdLineMapReduce.sh
2  ORIGINAL_FILE=$1
3  BLOCK_SIZE=$2
4  MapperParameters=$* #note $* will include the first and
second parameter also
5  CHUNK_FILE_PREFIX=$ORIGINAL_FILE.split
6  SORTED_CHUNK_FILES=$CHUNK_FILE_PREFIX*.sorted
7  usage ( )
8  {
9      echo Command Line Map Reduce
10     echo usage: cmdLineMapReduce filename chunksize
mapperParameters
11     echo Note: file1 will be split in chunks up to
$BLOCK_SIZE chunks each
12     echo NOTE: please make sure the output of the mapper is
what the reducer expects
13 }
14 #Splitting $ORIGINAL_FILE INTO CHUNKS
15 split -b $BLOCK_SIZE $ORIGINAL_FILE $CHUNK_FILE_PREFIX
16 #DISTRIBUTE
17 for file in $CHUNK_FILE_PREFIX*
18 do
19     #grep -i $FIND_WORD $file|wc -l >$file.intermediateCount
&
20     ./mapper.py $MapperParameters $file
>$file.intermediateCount &
21 done
22 wait
23 #MERGING INTERMEDIATE COUNT CAN TAKE THE FIRST COLUMN AND
TOTAL
24 #numOfInstances=$(cat *.intermediateCount | cut -f 1 | paste
&minussd+ &minus- |bc)
25 reducerOutput=$(cat *.intermediateCount | ./reducer.py)
26 echo "$reducerOutput"

```

Course Overview

- How to quantify data
- How to estimate time
- Powers of 10
- Parallelism
- Scale mundane tasks like grep, sort

Challenges in MapReduce Framework

- Data management
- Task management
- Fault tolerance
- Lack of interactivity

A Better Framework

- Scale out, not up
- Fault tolerance
- Minimize data movement
- Move processing to the data
- Hide file system details from application developer
- Seamless scalability

Class Phases

- Phase 0: Command line
- Phase 1: Hadoop/HDFS
- Phase 2: MRJob
- Phase 3: Spark

Summary

- Big data (what, why, where, who, how)
- Role of data scientist
- Data modeling pipeline
- Bias variance as a means of understanding more data
- MapReduce framework using command-line utility