

FIT5145 Introduction to Data Science

Module 2

Data Models in Organisations

2019 Lecture 4

Monash University

Student Feedback Survey

- ▶ Hope you enjoyed the unit so far!
- ▶ Spend a few mins now to fill in these two **anonymous** surveys
 - ▶ Lecture survey
 - ▶ Tutorial survey



Discussion: Python Language

- ▶ easy to learn
- ▶ flexible and multi-purpose
- ▶ great libraries
- ▶ well designed computer language
- ▶ good visualization for basic analysis

Reminder: Assessment

- ▶ all the Python you need to know is covered in tutorials
- ▶ **Assignment 1** will be the only assessment on Python code
 - ▶ but expect a question in exam like “compare Python as a data science tool with X” or “what features of Python make it good for Y”
- ▶ **the exam** is based on the lectures
 - ▶ use Alexandria as readings to support understanding
- ▶ **Assignments 2,4,5** are for a project proposal
 - ▶ you should be considering potential subject areas already
 - ▶ discuss them with tutor/lecturer

Unit Schedule: Modules

Module	Week	Content
1.	1	overview and look at projects (job) roles, and the impact
	2	
2.	3	data business models application areas and case studies
	4	
3.	5	characterising data and "big" data data sources and case studies
	6	
4.	7	resources and standards resources case studies
	8	
5.	10	data analysis theory data analysis process
	11	
6.	11	issues in data management GUEST SPEAKER & EXAM INFO
	12	

Learning Outcomes (Week 4)

By the end of this week you should be able to:

- ▶ Explain a set of criteria being used by NIST (National Institute of Standards and Technology) to analyse data science use cases
- ▶ Analyse a given use case based on a set of criteria used by NIST
- ▶ Use operations in R to read and manipulate data



Application Areas

(ePub section 2.5)

Consider different application areas:

- ▶ case studies from NIST:
 - ▶ provides a broad framework for analysing applications
- ▶ McKinsey Global Institute (MGI) report on big data:
 - ▶ study of different application areas

Application Areas: Case studies (from NIST)

provides a broad framework for analysing applications

NIST Case Studies

- ▶ this is the kind of analysis you will want in Assignment 2,4,5
- ▶ we will review their analysis

Caveat: many of the questions asked we won't look at properly until later modules

- ▶ so you may not be able to complete this kind of analysis now, but this should show you where we are headed
- ▶ in Assignment 2,4,5 you should be creative in adding more categories for your presentation

NIST Analysis

data sources: where does the data comes from?

data volume: how much there is?

data velocity: how does the data change over time?

data variety: what different kinds of data is there?

data veracity: is the data correct? what problems might it have?

software: what software needed to do the work?

analytics: what statistical analysis & visualisation is needed?

processing: what are the computational requirements?

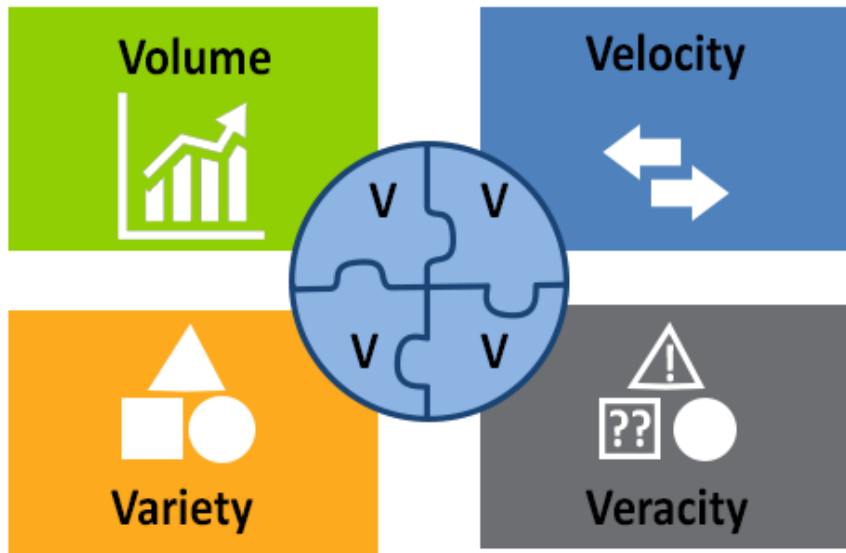
capabilities: what are key requirements of the operational system?

security/privacy: what security/privacy requirements are there?

lifecycle: what ongoing requirements are there?

other: are there other notable factors?

Four Vs of Big Data



NIST Analysis, cont.

Generally, we can relate the NIST categories to materials in our modules:

Module 2: capabilities; lifecycle; other

Module 3: data sources; data volume; data velocity; data variety; data veracity; processing;

Module 4: software; other

Module 5: analytics;

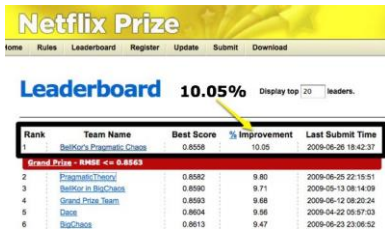
Module 6: security/privacy; lifecycle

Case Study: Netflix Movies

Netflix

On-demand internet streaming, and flat-rate DVD rental

- ▶ over 50 million subscribers in the US by 2014
- ▶ international market
- ▶ video recommendation!
- ▶ established the [Netflix Prize](#) in 2006-2009 as a crowdsourced way of testing out algorithms

A screenshot of the Netflix Prize Leaderboard. At the top, there's a yellow banner with 'Netflix Prize' in white. Below it is a navigation bar with links: Home, Rules, Leaderboard, Register, Update, Submit, Download. The 'Leaderboard' link is highlighted. Below the navigation bar, the word 'Leaderboard' is in large blue text, followed by '10.05%' in bold black text. To the right of '10.05%' is a small input field with '20' and the text 'Display top 20 leaders.' A yellow arrow points from the '10.05%' text to the '% Improvement' column header of the table below. The table has five columns: Rank, Team Name, Best Score, % Improvement, and Last Submit Time. The first row is highlighted in red and shows Rank 1, Team Name 'Belkor's Pragmatic Chaos', Best Score 0.8558, % Improvement 10.05, and Last Submit Time 2009-06-26 18:42:37. Below this row, a red banner reads 'Grand Prize - RMSE <= 0.8563'. The subsequent rows show Ranks 2 through 6 with their respective team names, best scores, improvement percentages, and last submit times.

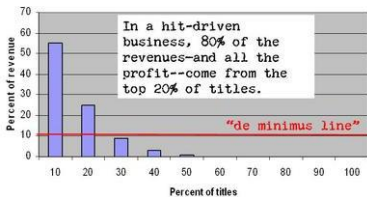
Rank	Team Name	Best Score	% Improvement	Last Submit Time
1	Belkor's Pragmatic Chaos	0.8558	10.05	2009-06-26 18:42:37
Grand Prize - RMSE <= 0.8563				
2	PragmaticTheory	0.8582	9.80	2009-06-25 22:15:51
3	Belkor in BigChaos	0.8590	9.71	2009-05-13 08:14:09
4	Grand Prize Team	0.8593	9.68	2009-06-12 08:20:24
5	Dane	0.8604	9.66	2009-04-22 05:57:03
6	BigChaos	0.8613	9.47	2009-06-23 23:06:52

By Ivongala (Own work) [Public domain], via Wikimedia Commons

Netflix: Background

Analysis follow the NIST Big Data WG Netflix analysis in [Volume 3, Use Cases and General Requirements](#), case 7 on page 8, A-24 and elsewhere

- ▶ Pareto principle, or **80/20 rule**:
 - ▶ top 20% of films watched 80% of time
 - ▶ standard video store stocked less than 20% of available titles in order to make the most money



from [The real meaning of 80/20](#)

- ▶ By adopting an Amazon style business model, Netflix could afford to rent the remaining 80%, the so-called long tail

FLUX Question



Based on the NIST analysis on page 8 of this week's slides.

Which of the following factor is least critical to Netflix application?

- A. Data volume
- B. Data velocity
- C. Data veracity
- D. Analytics

Netflix: Analysis

data sources: user movie ratings, user clicks, user profiles

data volume: in 2012: 25 million users, 4 million ratings/day, 3 million searches/day, video cloud storage of 2 petabytes

data velocity: video titles change daily, rankings/ratings updated

data variety: user rankings, user profiles, media properties

software: [Hadoop](#), [Pig](#), [Cassandra](#), [Teradata](#)

analytics: personalised recommender system

processing: analytic processing, streaming video

capabilities: ratings and search per day, content delivery

security/privacy: protect user data; digital rights

lifecycle: continued ranking and updating

other: mobile interface

Case Study: Electronic Medical Records (EMR)

EMR: Clinical Data



CLINICAL HISTORY: Cough, congestion.

COMMENTS:

PA and lateral views of chest reveals no evidence of active pleural or pulmonary parenchymal consolidation. There are diffusely increased interstitial lung markings consistent with chronic bronchitis. Underlying emphysema is not excluded. The cardiac silhouette is enlarged. The mediastinum and pulmonary vasculature are tortuous. Degenerative changes are noted in the thoracic spine.

IMPRESSION:

1. No evidence of acute pulmonary pathology.
2. Enlarged cardiac silhouette.
3. Tortuous aorta.
4. Diffusely increased interstitial lung markings consistent with chronic bronchitis. Underlying emphysema not excluded.
5. Consider follow up with Chest CT if clinically warranted.

EMR: Claims and Cost Data



P.O. BOX 742188
Los Angeles, CA 90074-2188

Customer Service
Toll Free: 1-800-549-3720, or
650-498-5850
9am-4pm
Monday-Friday

See Reverse Side for Patient Billing Details

Pay Online: stanfordhealthcare.org/billing
To activate your MyHealth account, visit
myhealth.stanfordhealthcare.org/activation and
enter the access code: XX0XX-0XXXX-XXX00

Guarantor ID# 000000000

YOUR PHYSICIAN STATEMENT

Page 1 of 1

Summary of Services and Amounts Due						
Service Date	Provider	Description	Charges	Credits	Insurance Balance	Patient Balance
Patient: SAMPLE, SAMPLE S		Visit #000000000	Service Line: General Surgery			
07/01/2014	Badger, James T, MD	99201 EVAL/MGMT OF NEW PATIENT	165.00		0.00	82.50
08/07/2014		UNINSURED DISCOUNT ADJ		-82.50		
		Totals:	165.00	-82.50	0.00	82.50

Electronic Medical Records

follows NIST Big Data WG Electronic Medical Records analysis in [Volume 3, Use Cases and General Requirements](#), case 16 in page 14, A-45 and elsewhere

Clinical data and claims/cost data is available per patient, per hospital

- ▶ large variety of sources of data
- ▶ systematic errors and difference in standards across institution

Task: segment patients into different types (“phenotypes”) to use in subsequent cohort studies

- ▶ case study is for Indiana Network for Patient Care

FLUX Question



Based on the NIST analysis on page 8 of this week's slides.

Which factor(s) is/are relevant to EMR case study?

- A. Data volume
- B. Data velocity
- C. Data variety
- D. None of the options
- E. All of the options

EMR: Analysis

data sources: clinical and claims data

data volume: 1000 centres, 12 million patients, 4 billion clinical events

data velocity: approx. 1 million clinical events/day

data variety: free text, lab results, pathology, outpatient, *etc.*

data veracity: different standards in different places

software: Hadoop, [Hive](#), Teradata, PostgreSQL, MongoDB

analytics: visualisation for data checking; standardisation of incoming data; general data analysis

processing: analytic processing, handling the volume

capabilities: models to support subsequent cohort studies

security/privacy: privacy and confidentiality required

lifecycle: full data management required

Case Study: Medical Imaging (MI)

MI Data



MI Task: Produce Analysis



CLINICAL HISTORY: Cough, congestion.

COMMENTS:

PA and lateral views of chest reveals no evidence of active pleural or pulmonary parenchymal consolidation. There are diffusely increased interstitial lung markings consistent with chronic bronchitis. Underlying pulmonary edema is not excluded. The cardiac silhouette is enlarged. The mediastinum and pulmonary vasculature are tortuous. Degenerative changes are noted in the thoracic spine.

IMPRESSION:

1. No evidence of acute pulmonary pathology.
2. Enlarged cardiac silhouette.
3. Tortuous aorta.
4. Diffusely increased interstitial lung markings consistent with chronic bronchitis. Underlying pulmonary edema is not excluded.
5. Consider follow up with Chest CT if clinically warranted.

Medical Imaging

follow NIST Big Data WG Pathology Imaging in [Volume 3, Use Cases and General Requirements](#), case 17 in page 14, A-48 and elsewhere

Biomedical data for imaging is high resolution and some is 3D

- ▶ interpretation of images done by trained experts
- ▶ requires significant training in interpretation
- ▶ many different kinds of instruments each requiring different interpretations
- ▶ millions produced daily in the USA

Medical Imaging: Analysis

data sources: biomedical image data

data volume: approx. 1 million events/day nationally

data variety: X-rays, CT scans, microscopes, ...

data veracity: current interpretation is often text based, so prone to text errors

software: advanced image processing and machine learning systems

analytics: computational image processing, supervised learning from images

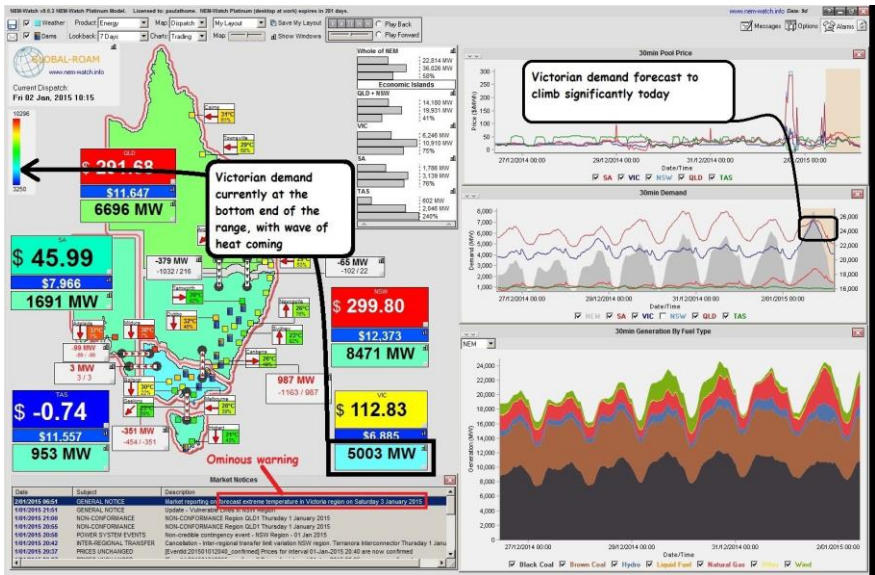
processing: handling the large volume, distributed and high throughput

capabilities: produce initial analysis for experts

security/privacy: privacy and confidentiality required

lifecycle: full data management required

Case Study: Electricity Demand Forecasting (EDF)



from [Analysing trends in VIC electricity demand](#)

Electricity Demand Forecasting

from NIST Big Data WG Electricity Demand Forecasting in

[Volume 3, Use Cases and General Requirements](#), case 51 in page 43 and A-134

Near realtime usage available thanks to smart meters

- ▶ with solar cells, consumers do energy generation too, but it is unpredictable
- ▶ main electricity generation must be planned
- ▶ brownouts and blackouts need to be prevented
- ▶ see see [Australian Energy Market Operator \(AEMO\)](#) and [their electricity site](#)

EDF: Analysis

data sources: utilities, smart meters, weather data, grids

data volume: city scale: 10GB/day

data velocity: updates every 15 minutes

data variety: time series, networks, spatial data

data veracity: occasional dropouts

software: advanced timeseries processing, spatial analysis

analytics: forecasting models

processing: handling the forecasting volume

capabilities: produce forecasts at different scales (hourly, daily)

security/privacy: privacy and confidentiality required

lifecycle: full data management required

Application Areas: McKinsey Global Institute report on big data

study of different application areas

Application Areas

We present details from McKinsey Global Institute report on Big Data from 2011, [*"Big data: The next frontier for innovation, competition, and productivity"*](#)

According to the MGI report, the main application areas of Big Data are:

1. Health
2. Government
3. Retail
4. Manufacturing
5. Location Technology

NB. What happened to Science? MGI is an industry organisation.

FLUX Question

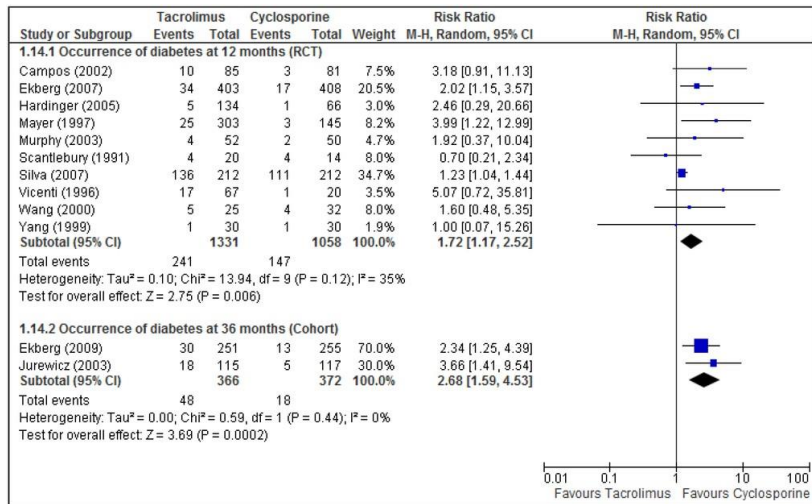


Give another example of Data Science applications that is not mentioned in this week's slides.

Application Areas: Health

Lets view [*"A Data Driven Approach to Diagnosing and Treating Disease"*](#), on [VideoLectures.NET](#) (video, see time 00:00-11.27)

Health: Pharmaceutical R&D Data



Health: Patient Behaviour and Sentiment Data



THE EXCHANGE BANK www.bocsonline.co
 Customer Service or Lost/Stolen: 800-693-1557
 Send Inquiries to: BANKERS CREDIT CARD SERVICE P.O. BOX 268856 OKLAHOMA CITY, OK 73126

ACCOUNT NUMBER	CREDIT LINE	AVAILABLE CREDIT	DAYS IN BILLING	STATEMENT CLOSING DATE	PAYMENT DUE DATE	MINIMUM PAYMENT DUE
	\$10,000	\$5,774.00	32	01/30/06	02/24/06	\$1,225.68

Tran Date	Post Date	Reference Number	Transaction Description	Amount
-----------	-----------	------------------	-------------------------	--------

PAYMENTS AND CREDITS

01/09	01/09	8543982QT06XBNV2B	PAYMENT - THANK YOU	1,458.09
-------	-------	-------------------	---------------------	----------

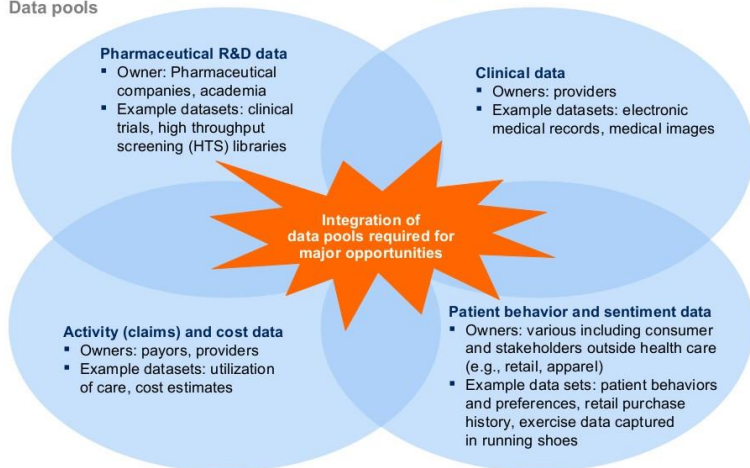
PURCHASES, DEBITS AND FINANCE CHARGES

01/05	01/05	8547082QR048Q1GFQ	MOGILLS TULSA OK	22.71
01/06	01/06	5542135QRWPGX31Z	MOLLY'S LANDING RESTAU GATOOSA OK	548.91
01/06	01/06	5548675QRZYX7P689	GOLDEN CORRAL FAMILY R OWASSO OK	9.7
01/11	01/11	5545370QW58HQ6007	CHELINO'S MEXOGAN REST OKLAHOMA CITY OK	29.7
01/12	01/12	0541019QX42QAEVEX	THE OLIVE GARD00015917 OKLAHOMA CITY OK	28.84
01/13	01/13	8548675QZ2YX7RP35	GOLDEN CORRAL FAMILY R OWASSO OK	9.7
01/13	01/13	5554186QZ03TM03TN	MARRIOTT HTL OKLAHOMA OKLAHOMA CITY OK	89.9
01/13	01/13	CHECK-IN 01/11/06	FOLIO #000005472	
01/13	01/13	5554186QZ03TM03TY	MARRIOTT HTL OKLAHOMA OKLAHOMA CITY OK	89.9
		CHECK-IN 01/11/06	FOLIO #000005472	

Health: Applications

Four distinct big data pools exist in the US health care domain today with little overlap in ownership and low integration

Data pools



SOURCE: McKinsey Global Institute analysis

Health: Clinical operations

Comparative effectiveness research: to study patient characteristics and the cost and outcomes of treatments

Clinical decision support systems: compare treatment against guidelines, to alert to drug interactions, *etc.*

Transparency about medical data: to help patients make more informed health care decisions, e.g., cesarean birth

Remote patient monitoring: support chronically ill, and at home care, to reduce subsequent hospital use

Advanced analytics applied to patient profiles: to identify who would benefit from proactive care or lifestyle changes

Health: Payment/Pricing

Automated systems: for insurance fraud detection and checking the accuracy and consistency of payors' claims

Health Economics: performance-based pricing plans based on real-world patient outcomes data to arrive at fair economic compensation

Health: Research and Development

Predictive modeling: rationale drug design

Statistical tools and algorithms to improve clinical trial design:
in the clinical phases of the R&D process

Analyzing clinical trials data: to identify additional indications
and discover adverse effects

Personalized medicine: understading genetic variation and

Analyzing disease patterns: analyzing disease patterns and
trends to model future demand/costs, and make
strategic R&D investment decisions

Application Areas: Government

lets review [*"The Mayor's Geek Squad"*](#) from New York Times

Government: Mayor's Geek Squad

- ▶ Lesson: analysis can be the basis for financial allocation decisions
- ▶ New York City was really the right city at the right time. Importantly that the mayor is a major advocate
- ▶ other cities:
 - ▶ infrastructures readiness, academic supports, citizen habits, and supporting mayor, cross-department cooperation
- ▶ data:
 - ▶ tollway, myki, traffic cameras, bike rental... social media, people "checking-in", or hashtag in tweets, facebook or instagram
- ▶ what is the nature of data availability?
 - ▶ city should provide enough public data to start from
 - ▶ help the local governing bodies plan urban development
 - ▶ could support international students better!

Government: Applications

What areas are the basis for good applications?

- ▶ Creating transparency
- ▶ Enabling experimentation to discover needs, expose variability, and improve performance
- ▶ Segmenting populations to customize actions
- ▶ Replacing/supporting human decision making with automated algorithms
- ▶ Innovating new business models, products, and services with big data

Application Areas: Retail

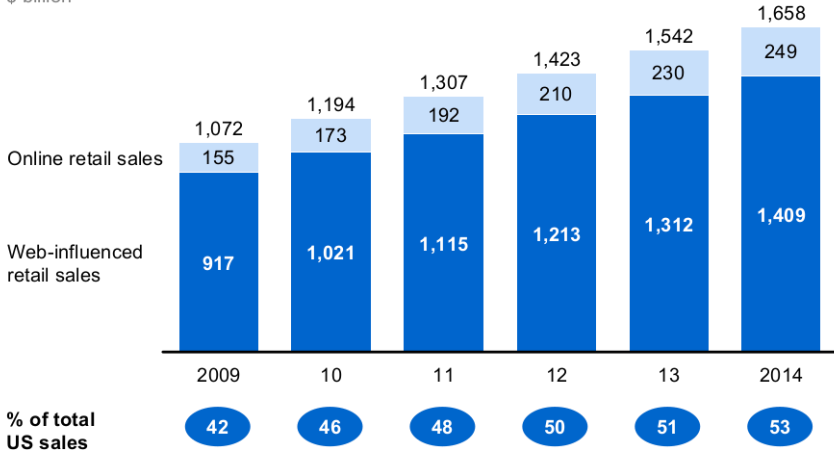
Understanding the Internet

- ▶ a major driver for new experiences in retail
- ▶ see ["How People Spend Their Time Online"](#) by GO-Gulf (infographic on a blog)
- ▶ see ["You are being watched"](#) by Techgenie, up on Pinterest.

Retail: Market Pressure

US online and Web-influenced retail sales are forecast to become more than half of all sales by 2013

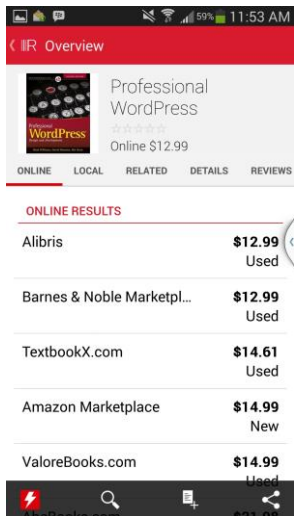
\$ billion



SOURCE: Forrester Research Web-influenced retail sales forecast, December 2009

Retail: Market Pressure, cont.

price comparison apps like this (from Ebay) let you be choosier



Retail: Applications

Marketing: Cross-selling; Location based marketing; In-store behavior analysis; Customer micro-segmentation; Sentiment analysis; Enhancing the multichannel; consumer experience

Merchandising: Assortment optimization; Pricing optimization; Placement and design; optimization

Operations: Performance transparency; Labor inputs optimization;

Supply chain: Inventory management; Distribution and logistics; optimization; Informing supplier negotiations

Application Areas: Manufacturing

Operations Analysis

The abundance and growth of machine data, which can include anything from IT machines to sensors and meters and GPS devices, is another major driver of big data solutions. In its raw format, many organizations are unable to leverage machine data. Yet disregarding this data means that organizations are making business decisions based only on a subset of available information. Leveraging machine data and combining it with existing enterprise data enables a new generation of applications that are able to analyze and gain insight from large volumes of multi-structured machine data—which in turn improves business results.

THE RESULTS

Empower the C-Suite

Reassure decision makers that they are acting with full knowledge & understanding of *all* available data.

Improve Reliability

Perform root cause analysis on data to more easily identify and preempt system failures, keeping customers happy.

Speed Operations

Help departments proactively minimize the problems and bottlenecks that stymie the flow of operations.

Monitor & React

Visualize streaming data to monitor the end-to-end infrastructure and deliver real-time alerts.

Raw Logs & Machine Data

Enterprise Data

Capture a Complete View

Access large volumes of machine, operational and transactional data and combine with other enterprise data.

Get the Context

Overcome complexities to perform advanced analysis and provide context across different data sets.

Get Insights From Analytics

Release intelligence trapped in your data, allowing agile interpretation and action.

WHAT DO YOU NEED TO SUCCEED?

Manufacturing Applications

We have identified the following big data levers across the manufacturing value chain

	R&D and design	Supply-chain mgmt	Production	Marketing and sales	After-sales service
1 Build consistent interoperable, cross-functional R&D and product design databases along supply chain to enable concurrent engineering, rapid experimentation and simulation, and co-creation	✓				
2 Aggregate customer data and make them widely available to improve service level, capture cross- and up-selling opportunities, and enable design-to-value	✓			✓	
3 Source and share data through virtual collaboration sites (idea marketplaces) to enable crowd sourcing)	✓			✓	
4 Implement advanced demand forecasting and supply planning across suppliers and using external variables		✓	✓	✓	
5 Implement lean manufacturing and model production virtually (digital factory) to create process transparency, develop dashboards, and visualize bottlenecks			✓		
6 Implement sensor data-driven operations analytics to improve throughput and enable mass customization			✓		
7 Collect after-sales data from sensors and feed back in real time to trigger after-sales services and detect manufacturing or design flaws			✓	✓	✓

SOURCE: McKinsey Global Institute analysis

Application Areas: Location

Mobile location-based services (LBS) and applications have proliferated

Mobile LBS applications continue to proliferate¹

People locating
(e.g., safety family/
child tracking,
friend finder)



Location check-in/
sharing on social
community
applications



City/regional guide,
neighborhood
service search



Location-enabled
entertainment, e.g.,
mobile gaming, geo-
tagged photo/travel



Revenue through the “Freemium” model

Revenue model for these mobile
LBS applications will be a mix of

- Free services/applications supported by advertising revenue
 - Sponsor links for mobile location-enabled (e.g., nearby point of interest) search
 - Advertising embedded in mobile applications
- Mobile apps requiring premiums for download or subscription
 - Onetime charge to download apps from mobile marketplaces
 - Recurring subscription fees for services/content
 - Add-on charges, e.g., purchase of virtual items in mobile games

¹ Navigation and other applications for non-individual usage have been assessed separately and are not included here.

How does location tracking work?

Input

RFID Tag



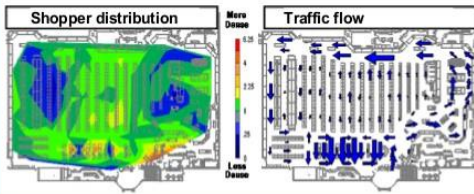
+ Mobile phones



Video

Personal tracking

Output



Location Applications

Location-based applications and services for individuals:

- ▶ smart routing
- ▶ automotive telematics
- ▶ mobile phone location-based services

Organizational use of individual personal location data:

- ▶ geo-targeted advertising
- ▶ electronic toll collection
- ▶ insurance pricing
- ▶ emergency response

Macro-level use of aggregate location data:

- ▶ urban planning.
- ▶ retail business intelligence
- ▶ some new business models

Next: Module 3

Data Types and Storage