# ภาษาการเขียนโปรแกรม 02-212-213

Week9

อ.ธิดาวรร คล้ายศรี



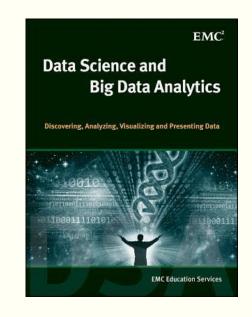


### Big Data Analytics

- Data Analytics Life Cycle
- Big data concepts & Hadoop ecosystem
- 3.1 Data Analytics Lifecycle (วงจรชีวิตการวิเคราะห์ข้อมูล)
  - 3.1.1 Describing Problem (กำหนดปัญหา)
  - 3.1.2 Identifying Data Sources (ระบุแหล่งข้อมูล)
  - 3.1.3 Preparing Data (การเตรียมความพร้อมข้อมูล)
  - 3.1.4 Planing Model (วางแผนโมเดล)
  - 3.1.5 Building Model (สร้างโมเดล)
  - 3.1.6 Visaulising Results (นำเสนอข้อมูลวิเคราะห์)
- 3.2 Big Data Concepts (แนวความคิดในเรื่อง big data)
- 3.3 Introductory Hadoop and Its Ecosystem (บทนำและวงชีวิตของฮาดูพ)

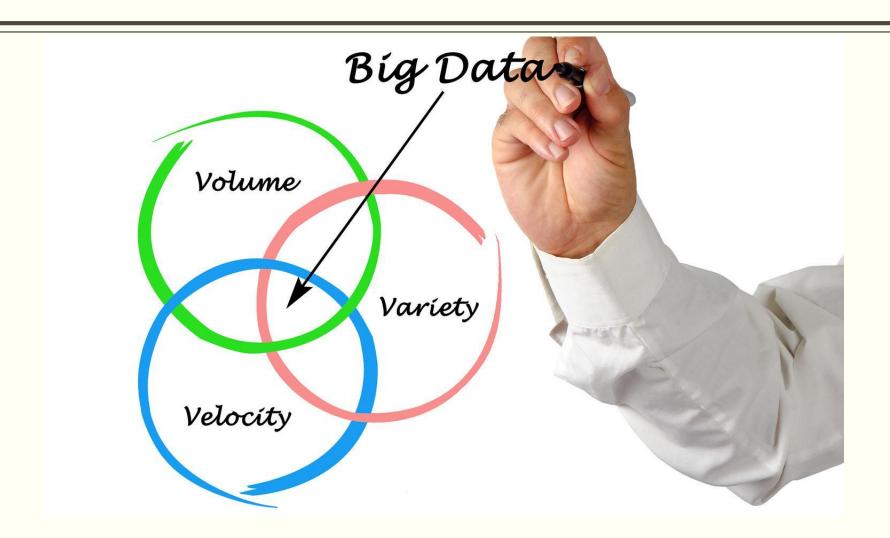
### Source of the content:

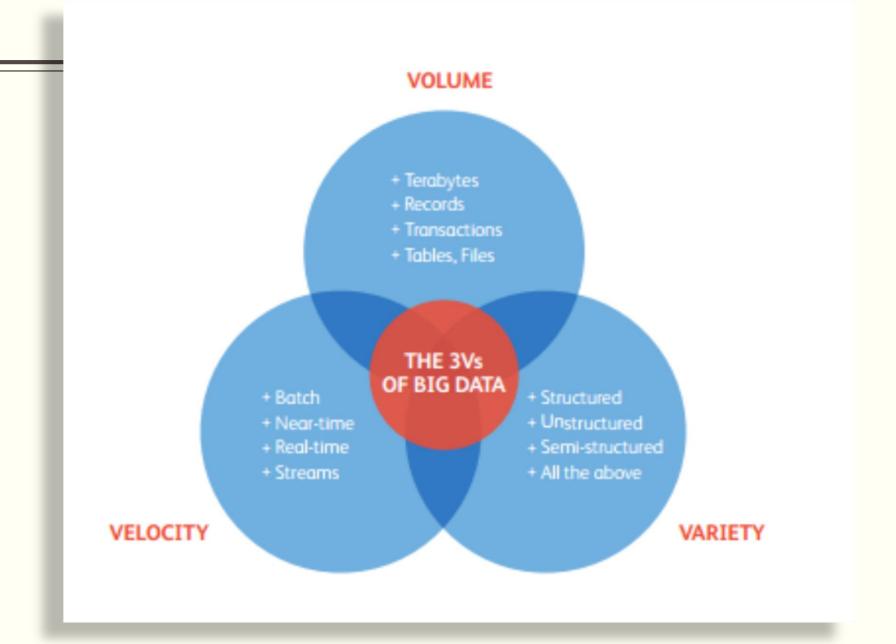
Majority of the slides: "Data Science and Big Data Analytic" text book by EMC Education Services A few slides: DEPA Big data analytic



Part 2: 3.2 Big Data Analytic Concepts

หลักในการวิเคราะห์ข้อมูลขนาดใหญ่





#### **40 ZETTABYTES** It's estimated that 2.5 QUINTILLION BYTES [ 43 TRILLION GIGABYTES ] [ 2.3 TRILLION GIGABYTES ] of data will be created by 2020, an increase of 300 of data are created each day times from 2005 Volume **6 BILLION** SCALE OF DATA have cell phones | Most companies in the U.S. have at least 100 TERABYTES 100,000 GIGABYTES I WORLD POPULATION: 7 BILLION of data stored Modern cars have close to The New York Stock Exchange 100 SENSORS captures 1 TB OF TRADE that monitor items such as fuel level and tire pressure INFORMATION during each trading session **Velocity** ANALYSIS OF STREAMING DATA

## The FOUR V's of Big **Data**

break big data into four dimensions: Volume, Velocity, Variety and Veracity

#### 4.4 MILLION IT JOBS



As of 2011, the global size of data in healthcare was estimated to be

#### 150 EXABYTES

[ 161 BILLION GIGABYTES ]



30 BILLION PIECES OF CONTENT are shared on Facebook every month

**Variety** 

# DIFFERENT

FORMS OF DATA

4 BILLION+ HOURS OF VIDEO

By 2014, it's anticipated

WEARABLE, WIRELESS

**HEALTH MONITORS** 

there will be

20 MILLION

are watched on YouTube each month



OO MILLION TWEETS

are sent per day by about 200 million monthly active users

### 1 IN 3 BUSINESS

don't trust the information they use to make decisions



in one survey were unsure of how much of their data was inaccurate



Poor data quality costs the US economy around

\$3.1 TRILLION A YEAR



Veracity UNCERTAINTY OF DATA

Sources: McKinsey Global Institute, Twitter, Cisco, Gartner, EMC, SAS, IBM, MEPTEC, QAS

By 2016, it is projected

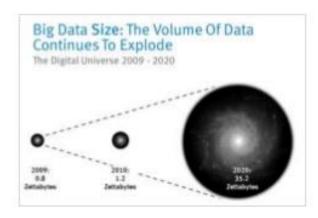
there will be

18.9 BILLION NETWORK CONNECTIONS - almost 2.5 connections per person on earth

### Key Characteristics of Big Data

#### Data Volume

 44x increase from 2010 to 2020 (1.2zettabytes to 35.2zb)



### 2. Processing Complexity

- Changing data structures
- Use cases warranting additional transformations and analytical techniques

#### 3. Data Structure

Greater variety of data structures to mine and analyze

# What's Driving The Data Deluge?



Mobile Sensors

READING SMART METERS
EVERY 15 MINUTES IS
3000X MORE
DATA INTENSIVE

Smart Grids

FACEBOOK UPLOADS 300 MILLION PHOTOS EACH DAY

Social Media

25000 DATA POINTS PER SECOND

Oil Exploration



Video Surveillance



Medical Imaging

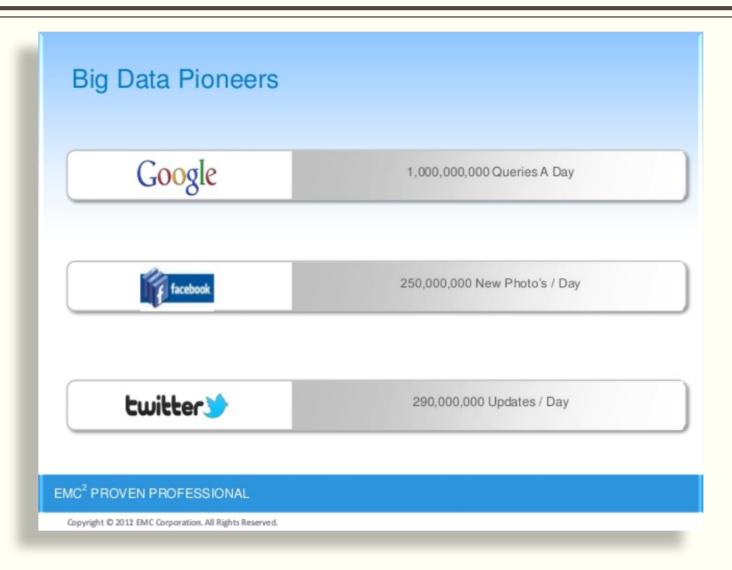


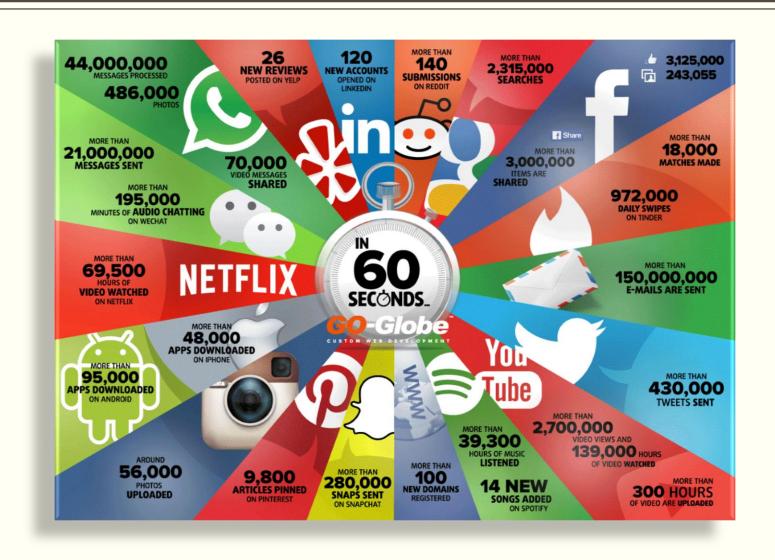
Video Rendering

ONE
GENOME
HAS FALLEN FROM
\$100M IN 2001
TO \$10K IN 2011

Gene Sequencing

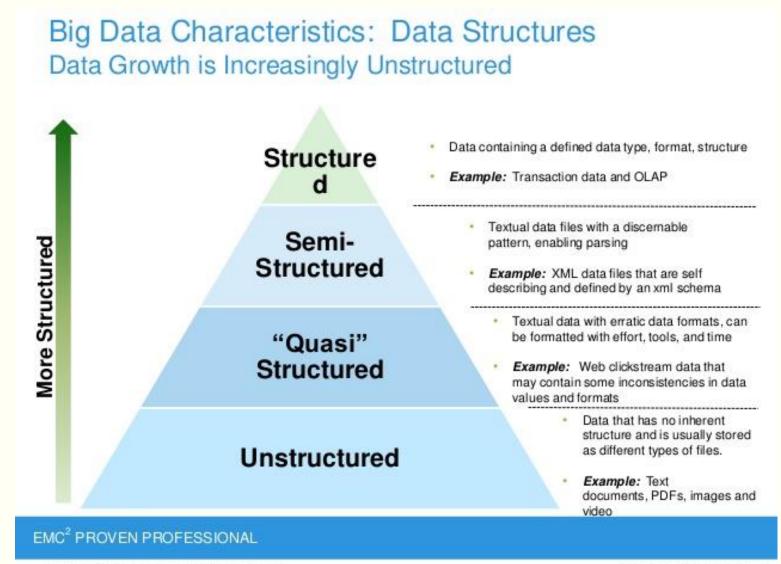
EMC<sup>2</sup>



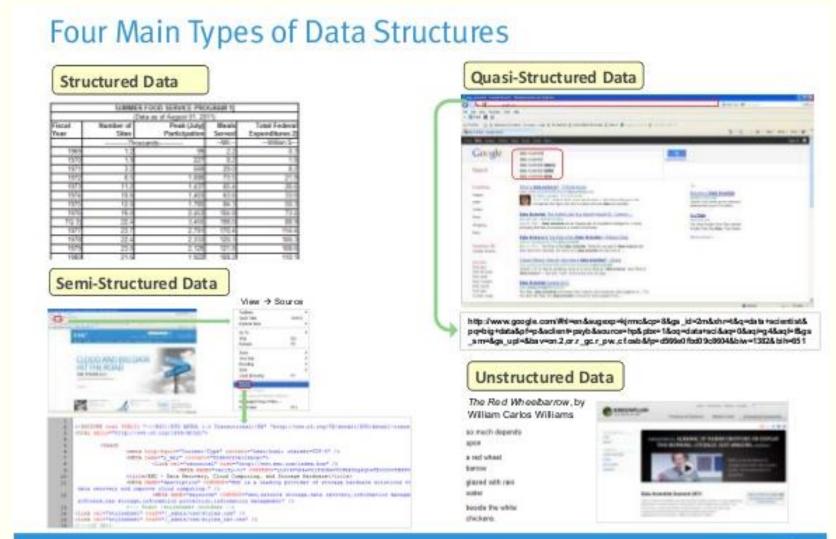




### 4 Types of Data (structures) in Big data era



### 4 Types of Data (structures) in Big data era





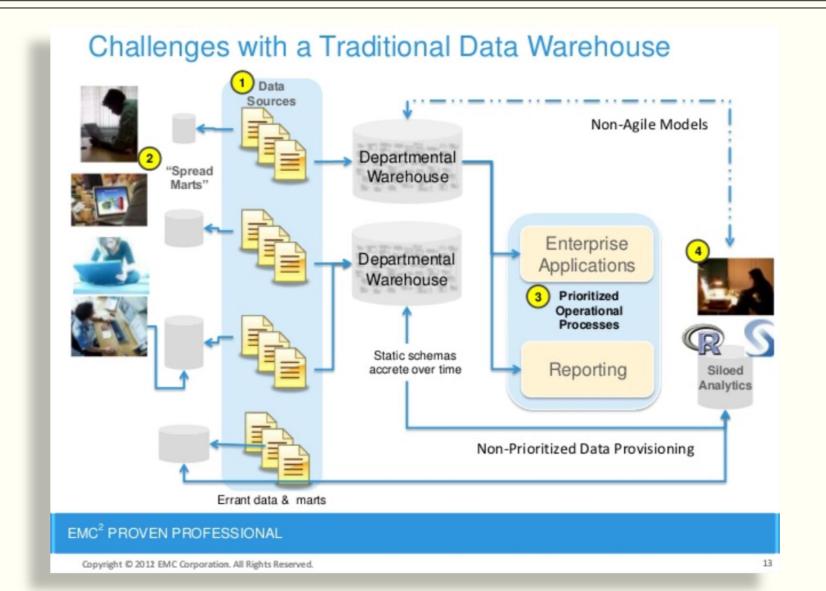
### Why-advanced data analytics?

### **Business Drivers for Advanced Analytics**

Current Business Problems Provide Opportunities for Organizations to Become More Analytical & Data Driven

_	Driver	Examples
1	Desire to optimize business operations	Sales, pricing, profitability, efficiency
2	Desire to identify business risk	Customer chum, fraud, default
3	Predict new business opportunities	Upsell, cross-sell, best new customer prospects
4	Comply with laws or regulatory requirements	Anti-Money Laundering, Fair Lending, Basel II

### Why-big data analytics?



### Why-big data analytics?

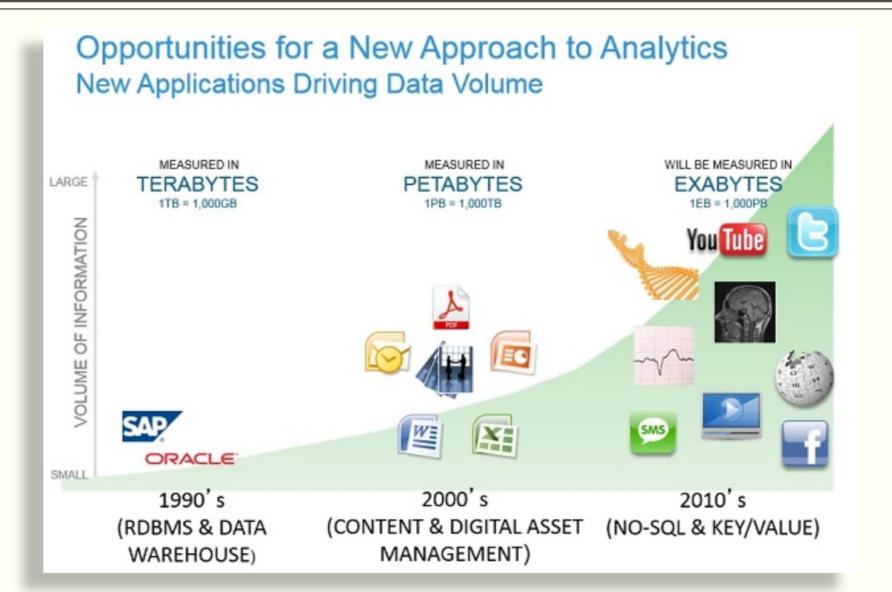
### Implications of a Traditional Data Warehouse

- High-value data is hard to reach and leverage
- Predictive analytics & data mining activities are last in line for data
  - Queued after prioritized operational processes
- Data is moving in batches from EDW to local analytical tools
  - In-memory analytics (such as R, SAS, SPSS, Excel)
  - Sampling can skew model accuracy
- Isolated, ad hoc analytic projects, rather than centrally-managed harnessing of analytics
  - Non-standardized initiatives
  - Frequently, not aligned with corporate business goals

Slow
"time-to-insight"

&
reduced
business impact

EMC<sup>2</sup> PROVEN PROFESSIONAL



### Why-big data analytics?

### Considerations for Big Data Analytics

### Criteria for Big Data Projects

- Speed of decision making
- Throughput
- Analysis flexibility

### **New Analytic Architecture**

#### **Analytic Sandbox**

Data assets gathered from multiple sources and technologies for analysis



- Enables high performance analytics using in-db processing
- Reduces costs associated with data replication into "shadow" file systems
- "Analyst-owned" rather than "DBA owned"

#### EMC<sup>2</sup> PROVEN PROFESSIONAL

### Big Data Analytics-use cases

### Big Data Analytics: Industry Examples

- Health Care
  - Reducing Cost of Care
- Public Services
  - Preventing Pandemics
- Life Sciences
  - Genomic Mapping
- IT Infrastructure
  - Unstructured Data Analysis
- Online Services
  - Social Media for Professionals



EMC2 PROVEN PROFESSIONAL



# Big Data Analytics: Healthcare



#### Situation

- Poor police response and problems with medical care, triggered by shooting of a Rutgers student
- The event drove local doctor to map crime data and examine local health care

**Use of Big Data** 

 Dr. Jeffrey Brenner generated his own crime maps from medical billing records of 3 hospitals

#### Key Outcomes

- City hospitals & ER's provided expensive care, low quality care
- Reduced hospital costs by 56% by realizing that 80% of city's medical costs came from 13% of its residents, mainly lowincome or elderly
- Now offers preventative care over the phone or through home visits

### EMC2 PROVEN PROFESSIONAL



# Big Data Analytics: Public Services



#### Situation

- Threat of global pandemics has increased exponentially
- Pandemics spreads at faster rates, more resistant to antibiotics

### Use of Big Data

- Created a network of viral listening posts
- Combines data from viral discovery in the field, research in disease hotspots, and social media trends
- Using Big Data to make accurate predications on spread of new pandemics
- · Identified a fifth form of human malaria, including its origin

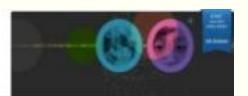
### Key Outcomes

- Identified why efforts failed to control swine flu
- Proposing more proactive approaches to preventing outbreaks

### EMC<sup>2</sup> PROVEN PROFESSIONAL



# Big Data Analytics: Life Sciences



#### Situation

Broad Institute (MIT & Harvard) mapping the Human Genome

### Use of Big Data

In 13 yrs, mapped 3 billion genetic base pairs; 8 petabytes

 Developed 30+ software packages, now shared publicly, along with the genomic data

### Key Outcomes

- Using genetic mappings to identify cellular mutations causing cancer and other serious diseases
- Innovating how genomic research informs new pharmaceutical drugs

#### EMC2 PROVEN PROFESSIONAL



# Big Data Analytics: IT Infrastructure



#### Situation

 Explosion of unstructured data required new technology to analyze quickly, and efficiently

### Use of Big Data

 Doug Cutting created Hadoop to divide large processing tasks into smaller tasks across many computers

#### Analyzes social media data generated by hundreds of thousands of users

#### Key Outcomes

- New York Times used Hadoop to transform its entire public archive, from 1851 to 1922, into 11 million PDF files in 24 hrs
- Applications range from social media, sentiment analysis, wartime chatter, natural language processing

### EMC2 PROVEN PROFESSIONAL



# Big Data Analytics: Online Services



Situation

Opportunity to create social media space for professionals

Use of Big Data

- Collects and analyzes data from over 100 million users
- Adding 1 million new users per week

Key Outcomes

- LinkedIn Skills, InMaps, Job Recommendations, Recruiting
- Established a diverse data scientist group, as founder believes this is the start of Big Data revolution

EMC<sup>2</sup> PROVEN PROFESSIONAL

# Case study: Starbuck-customer



# Case study: Real time-traffic



# Case study: Wholesale



# Case study: Hospital



# Case study: Banking

