Tamkang University



Big Data Mining

巨量資料探勘



巨量資料基礎: MapReduce典範、Hadoop與Spark生態系統 (Fundamental Big Data: MapReduce Paradigm, Hadoop and Spark Ecosystem)

1052DM02 MI4 (M2244) (3069) Thu, 8, 9 (15:10-17:00) (B130)



<u>Min-Yuh Day</u> <u>戴敏育</u> Assistant Professor 專任助理教授

Dept. of Information Management, Tamkang University

淡江大學 資訊管理學系



http://mail.tku.edu.tw/myday/ 2017-02-23

課程大綱 (Syllabus)

週次(Week) 日期(Date) 內容(Subject/Topics)

- 1 2017/02/16 巨量資料探勘課程介紹 (Course Orientation for Big Data Mining)
- 2 2017/02/23 巨量資料基礎: MapReduce典範、Hadoop與Spark生態系統 (Fundamental Big Data: MapReduce Paradigm, Hadoop and Spark Ecosystem)
- 3 2017/03/02 關連分析 (Association Analysis)
- 4 2017/03/09 分類與預測 (Classification and Prediction)
- 5 2017/03/16 分群分析 (Cluster Analysis)
- 6 2017/03/23 個案分析與實作一 (SAS EM 分群分析): Case Study 1 (Cluster Analysis – K-Means using SAS EM)
- 7 2017/03/30 個案分析與實作二 (SAS EM 關連分析): Case Study 2 (Association Analysis using SAS EM)

課程大綱 (Syllabus)

週次(Week) 日期(Date) 內容(Subject/Topics)

- 8 2017/04/06 教學行政觀摩日 (Off-campus study)
- 9 2017/04/13 期中報告 (Midterm Project Presentation)
- 10 2017/04/20 期中考試週 (Midterm Exam)
- 11 2017/04/27 個案分析與實作三 (SAS EM 決策樹、模型評估): Case Study 3 (Decision Tree, Model Evaluation using SAS EM)
- 12 2017/05/04 個案分析與實作四 (SAS EM 迴歸分析、類神經網路): Case Study 4 (Regression Analysis, Artificial Neural Network using SAS EM)
- 13 2017/05/11 Google TensorFlow 深度學習 (Deep Learning with Google TensorFlow)
- 14 2017/05/18 期末報告 (Final Project Presentation)
- 15 2017/05/25 畢業班考試 (Final Exam)

2017/02/23 巨量資料基礎: MapReduce典範、 Hadoop與Spark生態系統 (Fundamental Big Data: **MapReduce** Paradigm, Hadoop and Spark Ecosystem

Big Data Analytics and **Data Mining**

Architectures of Big Data Analytics

Architecture of Big Data Analytics



Architecture of Big Data Analytics



Architecture for Social Big Data Mining

(Hiroshi Ishikawa, 2015)



Business Intelligence (BI) Infrastructure



Data Warehouse Data Mining and Business Intelligence



The Evolution of BI Capabilities



Source: Turban et al. (2011), Decision Support and Business Intelligence Systems

Data Science and Business Intelligence



Data Science and Business Intelligence



Predictive Analytics and Data Mining (Data Science)

Time

Future

Past

Predictive Analytics and Data Mining (Data Science)

Structured/unstructured data, many types of sources, very large datasets

Optimization, predictive modeling, forecasting statistical analysis

What if...?

What's the optimal scenario for our business? What will happen next? What if these trends countinue? Why is this happening?

Data Mining at the Intersection of Many Disciplines



Source: Turban et al. (2011), Decision Support and Business Intelligence Systems

A Taxonomy for Data Mining Tasks



Source: Turban et al. (2011), Decision Support and Business Intelligence Systems

Traditional Analytics



Hadoop as a "new data" Store



Hadoop as an additional input to the EDW



Source: Deepak Ramanathan (2014), SAS Modernization architectures - Big Data Analytics

Hadoop Data Platform As a "staging Layer" as part of a "data Lake"

- Downstream stores could be Hadoop, data appliances or an RDBMS



SAS Big data Strategy - SAS areas



Source: Deepak Ramanathan (2014), SAS Modernization architectures - Big Data Analytics

SAS Big data Strategy - SAS areas



Source: Deepak Ramanathan (2014), SAS Modernization architectures - Big Data Analytics

SAS[®] Within the HADOOP ECOSYSTEM



SAS enables the entire lifecycle around HADOOP

SAS enableS the entire lifecycle around HADOOP

Done using either the Data Preparation, Data Exploration or Build Model Tools SAS Visual Analytics **Decision Manager** PROBLEM EVALUATE / MONITOR RESULTS SAS Visual Analytics DATA PREPARATION SAS Visual Statistics SAS In-Memory Statistics for Hadoop SAS Scoring Accelerator for Hadoop SAS Code Accelerator for Hadoop Done using either the Data Preparation, Data Exploration or Build Model Tools **Decision Manager** SAS High Performance Analytics Offerings supported by relevant clients like SAS Enterprise Miner, SAS/STAT etc.

Data Mining Process

Data Mining Process

- A manifestation of best practices
- A systematic way to conduct DM projects
- Different groups has different versions
- Most common standard processes:
 - CRISP-DM
 - (Cross-Industry Standard Process for Data Mining)
 - SEMMA
 - (Sample, Explore, Modify, Model, and Assess)
 - KDD

(Knowledge Discovery in Databases)

Data Mining Process (SOP of DM) What main methodology are you using for your analytics, data mining, or data science projects ?

Data Mining Process

CRISP-DM (86)	43% 42%
My own (55)	27.5%
SEMMA (17)	8.5% 13%
Other, not domain-specific (16)	8% 4%
KDD Process (15)	7.5% 7.3%
My organizations' (7)	3.5% 5.3%
A domain-specific methodology (4)	2% 4.7%
None (0)	0% 4.7%

2014 poli 2007 poli





Data Mining: Core Analytics Process

The KDD Process for Extracting Useful Knowledge from Volumes of Data

Source: Fayyad, U., Piatetsky-Shapiro, G., & Smyth, P. (1996). The KDD Process for Extracting Useful Knowledge from Volumes of Data. Communications of the ACM, 39(11), 27-34.

Fayyad, U., Piatetsky-Shapiro, G., & Smyth, P. (1996). The KDD Process for **Extracting Useful Knowledge** from Volumes of Data. Communications of the ACM, 39(11), 27-34.

Knowledge Discovery in Databases creates the context for developing the tools needed to control the flood of data facing organizations that depend on ever-growing databases of business, manufacturing, scientific, and personal information.

The KDD Process for Extracting Useful Knowledge from Volumes of Data

As we march into the age of digital information, the problem of data overload looms ominously ahead. Our ability to analyze and Gregory Piatetsky-Shapiro, understand massive datasets lags far behind our ability to gather and store the data. A new gen-

the rapidly growing volumes of data. data warehouses. These techniques and tools are the Current hardware and database techdata mining

and Padhraic Smyth

Usama Fayyad,

eration of computational techniques and many more applications generate and tools is required to support the streams of digital records archived in extraction of useful knowledge from huge databases, sometimes in so-called

subject of the emerging field of knowl- nology allow efficient and inexpensive edge discovery in databases (KDD) and reliable data storage and access. However er, whether the context is business Large databases of digital informa- medicine, science, or government, the tion are ubiquitous. Data from the datasets themselves (in raw form) are of neighborhood store's checkout regis- little direct value. What is of value is the ter, your bank's credit card authoriza- knowledge that can be inferred from tion device, records in your doctor's the data and put to use. For example, office, patterns in your telephone calls, the marketing database of a consumer

Data Mining

Knowledge Discovery in Databases (KDD) Process

(Fayyad et al., 1996)



Knowledge Discovery in Databases (KDD) Process



Source: Jiawei Han and Micheline Kamber (2006), Data Mining: Concepts and Techniques, Second Edition, Elsevier

Data Mining Process: CRISP-DM



Source: Turban et al. (2011), Decision Support and Business Intelligence Systems

Data Mining Process: CRISP-DM

- **Step 1:** Business Understanding
- Step 2: Data Understanding
- Step 3: Data Preparation (!)
- Step 4: Model Building
- **Step 5:** Testing and Evaluation
- Step 6: Deployment
- The process is highly repetitive and experimental (DM: art versus science?)



Data Preparation – A Critical DM Task



Source: Turban et al. (2011), Decision Support and Business Intelligence Systems


Source: Turban et al. (2011), Decision Support and Business Intelligence Systems

Data Mining Processing Pipeline

(Charu Aggarwal, 2015)



Source: Charu Aggarwal (2015), Data Mining: The Textbook Hardcover, Springer

Fundamental Big Data: MapReduce Paradigm, Hadoop and Spark **Ecosystem**



MapReduce Paradigm

MapReduce Paradigm



MapReduce Word Count

Input

Dog Love Cat Bird Love Bird Dog Bird Cat

MapReduce Word Count

Input

Dog Love Cat

Bird Love Bird

Dog Bird Cat

Bird, 3 Cat, 2 Dog, 2 Love, 2

Output

MapReduce Word Count



Hadoop Ecosystem



The Apache[™] Hadoop[®] project develops open-source software for reliable, scalable, distributed computing.





Big Data with Hadoop Architecture

LOGICAL ARCHITECTURE





PHYSICAL ARCHITECTURE





Hadoop Cluster

Big Data with Hadoop Architecture Logical Architecture Processing: MapReduce



Big Data with Hadoop Architecture Logical Architecture Storage: HDFS



Big Data with Hadoop Architecture Process Flow



Big Data with Hadoop Architecture Hadoop Cluster



Hadoop Ecosystem





HDP (Hortonworks Data Platform) A Complete Enterprise Hadoop Data Platform



Apache Hadoop Hortonworks Data Platform



Hadoop and Data Analytics Tools



Hadoop 1 \rightarrow Hadoop 2

Hadoop 1

- Silos & Largely batch
- Single Processing engine

Script	SQL	Real-time	Others
Pig	Hive	HBase	Storm, Solr, etc.
MapReduce (Cluster Resource Management & Data Processing)			
1			
HDFS			
(Hadoop Distributed File System)			

Hadoop 2 w/

- Multiple Engines, Single Data Set
- · Batch, Interactive & Real-Time



Big Data Solution



Source: http://www.newera-technologies.com/big-data-solution.html

Traditional ETL Architecture



Offload ETL with Hadoop (Big Data Architecture)



Spark Ecosystem



Lightning-fast cluster computing

Apache Spark is a fast and general engine for large-scale data processing.



Logistic regression in Hadoop and Spark



Run programs up to 100x faster than Hadoop MapReduce in memory, or 10x faster on disk.

Source: http://spark.apache.org/



Ease of Use

• Write applications quickly in Java, Scala, Python, R.



text_file = spark.textFile("hdfs://...")

text_file.flatMap(lambda line: line.split())
.map(lambda word: (word, 1))
.reduceByKey(lambda a, b: a+b)

Spark and Hadoop













Spark Ecosystem

Spark Spark SQL Streaming

MLlib (machine learning)

GraphX (graph)

Apache Spark



Source: Mike Frampton (2015), Mastering Apache Spark, Packt Publishing

Hadoop vs. Spark



Steps to **Install Hadoop** on a **Personal Computer** (Windows/OS X)

Hodoop: Linux Based Software






Appliance



Connection to Hadoop





Source: https://www.youtube.com/watch?v=rO-V1mxhzcM&list=PLyZEf-TOnZen8E5m5TlpIsdok2fyKDNRa&index=5

Virtual Box



About

Screenshots

Downloads

Contribute

Community

Documentation

End-user docs

Technical docs



Welcome to VirtualBox.org!

VirtualBox is a powerful x86 and AMD64/Intel64 virtualization product for enterprise as well as home use. Not only is VirtualBox an extremely feature rich, high performance product for enterprise customers, it is also the only professional solution that is freely available as Open Source Software under the terms of the GNU General Public License (GPL) version 2. See "About VirtualBox" for an introduction.

Presently, VirtualBox runs on Windows, Linux, Macintosh, and Solaris hosts and supports a large number of guest operating systems including but not limited to Windows (NT 4.0, 2000, XP, Server 2003, Vista, Windows 7, Windows 8, Windows 10), DOS/Windows 3.x, Linux (2.4, 2.6, 3.x and 4.x), Solaris and OpenSolaris, OS/2, and OpenBSD.

VirtualBox is being actively developed with frequent releases and has an ever growing list of features, supported guest operating systems and platforms it runs on. VirtualBox is a community effort backed by a dedicated company: everyone is encouraged to contribute while Oracle ensures the product always meets professional quality criteria.



Hot picks:

- Pre-built virtual machines for developers at
 → Oracle Tech Network
- Hyperbox Open-source Virtual Infrastructure Manager ⇒ project site
- **phpVirtualBox** AJAX web interface ⇒project site
- IQEmu automated Windows VM creation, application integration ⇒http://mirage335-site.member.hacdc.org:6380/wiki/Category:IQEmu



https://www.virtualbox.org/

search... Login Preferences

News Flash

- New January 17th, 2017 VirtualBox 5.1.14 released! Oracle today released a 5.1 maintenance release which improves stability and fixes regressions. See the Changelog for details.
- Important December 2nd, 2016
 We're hiring!
 Looking for a new challenge? We're looking for a GUI developer
 (Germany/European Union).
- New July 12th, 2016 VirtualBox 5.1 released! Many enhancements and improvements. Read more in the announcement.

More information...



Source: https://www.youtube.com/watch?v=rO-V1mxhzcM&list=PLyZEf-TOnZen8E5m5TlpIsdok2fyKDNRa&index=5

Hortonworks Sandbox

The easiest way to get started with Enterprise Hadoop



http://hortonworks.com/products/hortonworks-sandbox/#install

Get started on Hadoop with these tutorials based on the Hortonworks Sandbox

	COMMUNITY	BLOGS	PARTNERS	CONTACT US	Q	SUPPORT LOGIN	E	NGLISH 🛩
HORTONWORKS	Products	s Soluti	ons Cu	istomers	Services & Sup	port About	Us	GET STARTED

TUTORIALS

Get started on Hadoop with these tutorials based on the Hortonworks Sandbox



DEVELOP WITH HADOOP

Start developing with Hadoop. These tutorials are designed to ease your way into developing with Hadoop:

Apache Spark on HDP

1

Hands-on Tour of Apache Spark in 5 Minutes

Apache Spark is a fast, in-memory data processing engine with elegant and expressive development APIs in Scala, Java, Python, and R that allow data workers to efficiently execute machine learning algorithms that require fast iterative access to datasets (see Spark API Documentation for more info). Spark on Apache Hadoop YARN enables deep integration with Hadoop [...]

SHARE in f 🎔 🛛 NEWSLETTER 💟

Contact Sales?

Apache Hadoop

Apache > Hadoop > Search with Apache Solr Search Wiki Top Last Published: 01/27/2017 02:33:46 About Welcome What Is Apache Welcome to Apache[™] Hadoop[®]! Hado... Getting Started ... PDF Download Hadoop Who Uses Hadoop?... What Is Apache Hadoop? News Releases The Apache[™] Hadoop® project develops open-source software for reliable, scalable, distributed computing. **Release Versioning** The Apache Hadoop software library is a framework that allows for the distributed processing of large data sets across clusters of computers using simple Mailing Lists programming models. It is designed to scale up from single servers to thousands of machines, each offering local computation and storage. Rather than rely on Issue Tracking hardware to deliver high-availability, the library itself is designed to detect and handle failures at the application layer, so delivering a highly-available service on Who We Are? top of a cluster of computers, each of which may be prone to failures. Who Uses Hadoop? Buy Stuff The project includes these modules: Sponsorship Thanks Hadoop Common: The common utilities that support the other Hadoop modules. Privacy Policy Hadoop Distributed File System (HDFS™): A distributed file system that provides high-throughput access to application data. Bylaws Hadoop YARN: A framework for job scheduling and cluster resource management. Committer criteria . License Hadoop MapReduce: A YARN-based system for parallel processing of large data sets. . Documentation Other Hadoop-related projects at Apache include: **Related Projects** Ambari™: A web-based tool for provisioning, managing, and monitoring Apache Hadoop clusters which includes support for Hadoop HDFS, Hadoop built with MapReduce, Hive, HCatalog, HBase, ZooKeeper, Oozie, Pig and Sqoop. Ambari also provides a dashboard for viewing cluster health such as heatmaps and Apache Forrest ability to view MapReduce, Pig and Hive applications visually alongwith features to diagnose their performance characteristics in a user-friendly manner. • Avro[™]: A data serialization system. **Cassandra**^M: A scalable multi-master database with no single points of failure. ٠ Chukwa[™]: A data collection system for managing large distributed systems.

- **HBase**[™]: A scalable, distributed database that supports structured data storage for large tables.
- Hive^m: A data warehouse infrastructure that provides data summarization and ad hoc querying.
- Mahout[™]: A Scalable machine learning and data mining library.
- **<u>Pig</u>**[™]: A high-level data-flow language and execution framework for parallel computation.

http://hadoop.apache.org/

Apache Hadoop

http://hadoop.apache.org/releases.html#Download



- Committer criteria
- License
- Documentation
- Related Projects

version	Release Date	larball	GPG	SHA-256
3.0.0-alpha2	25 January, 2017	<u>source</u>	<u>signature</u>	checksum file
		<u>binary</u>	<u>signature</u>	checksum file
3.0.0-alpha1	03 September, 2016	<u>source</u>	<u>signature</u>	checksum file
		<u>binary</u>	<u>signature</u>	checksum file
<u>2.7.3</u>	25 August, 2016	source	<u>signature</u>	227785DC 6E3E6EF8
		<u>binary</u>	<u>signature</u>	D489DF38 08244B90
2.6.5	08 October, 2016	source	<u>signature</u>	3A843F18 73D9951A
		<u>binary</u>	signature	001AD18D 4B6D0FE5
2.5.2	19 Nov, 2014	source	<u>signature</u>	139EF872 09C5637E
		binary	signature	0BDB4850 A3825208

To verify Hadoop releases using GPG:

- 1. Download the release hadoop-X.Y.Z-src.tar.gz from a mirror site.
- 2. Download the signature file hadoop-X.Y.Z-src.tar.gz.asc from Apache.
- 3. Download the Hadoop KEYS file.
- 4. gpg --import KEYS
- 5. gpg --verify hadoop-X.Y.Z-src.tar.gz.asc

To perform a quick check using SHA-256:

- 1. Download the release hadoop-X.Y.Z-src.tar.gz from a mirror site.
- 2. Download the checksum hadoop-X.Y.Z-src.tar.gz.mds from Apache.

Apache Hadoop YARN



Source: http://hadoop.apache.org/docs/current/hadoop-yarn/hadoop-yarn-site/YARN.html

Apache Spark



Ease of Use

data flow and in-memory computing.

Write applications quickly in Java, Scala, Python, R.

Apache Spark has an advanced DAG execution engine that supports acyclic

Spark offers over 80 high-level operators that make it easy to build parallel apps. And you can use it *interactively* from the Scala, Python and R shells.

0.9 0

Logistic regression in Hadoop and Spark

text_file = spark.textFile("hdfs://...")

text_file.flatMap(lambda line: line.split()) .map(lambda word: (word, 1)) .reduceByKey(lambda a, b: a+b)

Word count in Spark's Python API

http://spark.apache.org/

Apache Software Foundation -

Latest News

Spark Summit East (Feb 7-9th, 2017, Boston) agenda posted (Jan 04, 2017)

Spark 2.1.0 released (Dec 28, 2016)

Spark wins CloudSort Benchmark as the most efficient engine (Nov 15, 2016)

Spark 2.0.2 released (Nov 14, 2016)

Archive

Download Spark

Built-in Libraries:

SQL and DataFrames Spark Streaming MLlib (machine learning) GraphX (graph)

Third-Party Projects

References

- EMC Education Services (2015), Data Science and Big Data Analytics: Discovering, Analyzing, Visualizing and Presenting Data, Wiley
- Shiva Achari (2015), Hadoop Essentials - Tackling the Challenges of Big Data with Hadoop, Packt Publishing
- Mike Frampton (2015), Mastering Apache Spark, Packt Publishing
- Deepak Ramanathan (2014), SAS Modernization architectures - Big Data Analytics, http://www.slideshare.net/deepakramanathan/sasmodernization-architectures-big-data-analytics