

# COMM1822

Term 2 2022

## Introduction to Databases for Business Analytics

Assignment Project Exam Help

Week 9 Big Data 2

<https://eduassistpro.github.io/>

Add WeChat edu\_assist\_pro

Lecturer-in-Charge: Kam-Fung (Henry) Cheung

Email: [kf.cheung@unsw.edu.au](mailto:kf.cheung@unsw.edu.au)

Tutors: Theresa Tran

Liam Li Chen

Kathy Xu

PASS Leader: Srilekha Chandrashekara Kolaki



## Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu\_assist\_pro

# Copyright

- There are some file-sharing websites that specialise in buying and selling academic work to and from university students.

**Assignment Project Exam Help**

- If you upload your original work and presents it as their own either on a file-sharing website or on a social media platform, you may be found guilty of collusion — even years after graduation

**Add WeChat edu\_assist\_pro**

- These file-sharing websites may also accept purchase of course materials, **such as copies of lecture slides and tutorial handouts**. By law, the copyright on course materials, developed by UNSW staff in the course of their employment, belongs to UNSW. It constitutes copyright infringement, if not academic misconduct, to trade these materials.

# Acknowledgement of Country

UNSW Business School acknowledges the Bidjigal (Kensington campus) and Gadigal (City campus) the traditional custodians of the lands where each campus is located.

Assignment Project Exam Help

We acknowledge all Aboriginal and Torres Strait Islander Elders, past and present and their communities who have shared and practiced their teachings over thousands of years including business practices.

<https://eduassistpro.github.io/>

Add WeChat edu\_assist\_pro

We recognise Aboriginal and Torres Strait Islander people's ongoing leadership and contributions, including to business, education and industry.

UNSW Business School. (2022, May 7). *Acknowledgement of Country* [online video]. Retrieved from <https://vimeo.com/369229957/d995d8087f>

## Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu\_assist\_pro

Chapter 14  
Assignment Project Exam Help  
and NoSQL  
<https://eduassistpro.github.io/>  
Add WeChat edu\_assist\_pro

# W9 Learning Outcomes

## ☐ **Big Data Technologies**

### ☐ Hadoop Ecosystem

- ☐ Hadoop Distributed File System (HDFS)
- ☐ MapReduce
- ☐ Pig
- ☐ Hive
- ☐ HBase
- ☐ Impala

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu\_assist\_pro

## ☐ **NoSQL Database Types**

- ☐ Key-value databases
- ☐ Document databases
- ☐ Column-oriented databases
- ☐ Graph databases

## ☐ **Big Data Strategies**

# Big Data Technologies

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu\_assist\_pro



**UNSW**  
SYDNEY



# Big Data Infrastructure Challenges

## ❑ Linear scalability

- ❑ To accommodate for the scalability of processing, thereby the storage management and architecture of traditional data management techniques become obsolete.

## ❑ High throughput

- ❑ Infrastructure that is extremely processing, and storage.

## ❑ Fault tolerance

- ❑ Any portion of the processing architecture should be able to take over and resume processing from the point of failure in any other part of the system.

# Big Data Infrastructure Challenges

## ❑ Auto recovery

- ❑ The processing architecture should be self-managing and recover from failure without manual intervention.

Assignment Project Exam Help

## ❑ High degree of parallelism

<https://eduassistpro.github.io/>

- ❑ Distribute the load across multiple machines, each processing a different program. e.g., data analysis uods: linear regression, random forests

Add WeChat edu\_assist\_pro

## ❑ Distributed data processing

- ❑ The underlying platform must be able to process distributed data to achieve extreme scalability.

# What is Hadoop?



- ❑ Hadoop is an open-source framework for storing and analyzing massive amounts of distributed, **unstructured** data.

**Assignment Project Exam Help**

- ❑ Hadoop was created by D. **arella** in 2005.

**<https://eduassistpro.github.io/>**

- ❑ Hadoop clusters run on inexpensive computers so projects can scale-out inexpensively.

**Add WeChat edu\_assist\_pro**

- ❑ Open source - hundreds of contributors continuously improve the core technology.

- ❑ What is Hadoop? - <https://www.youtube.com/watch?v=9s-vSeWej1U>

# Hadoop

- ❑ Not a single product, not a single database.
- ❑ A **collection of big data applications**.
- ❑ A **framework**, platform and ecosystem.
- ❑ Consisting of **different** <https://eduassistpro.github.io/>
- ❑ Most important components:
  - Hadoop Distributed File System (HDFS)
  - MapReduce
  - Pig
  - Hive
  - HBase
  - Impala

Assignment Project Exam Help

Add WeChat edu\_assist\_pro

# Why Hadoop?

## ❑ Problems with relational database management system (RDBMS):

- Insufficiently scalable for big data
  - Insufficient speed for live data
  - Lack of sophisticated
  - Essentially a design b
- can easily “scale up” to an extend, but scale out”)
- Assignment Project Exam Help  
<https://eduassistpro.github.io/>  
Add WeChat edu\_assist\_pro
- CPU and RAM (you

## ❑ Polyglot persistence: The coexistence of a variety of data storage and data management technologies within an organization's infrastructure.

Structured: Customer's data, e.g., date of birth, address, bank account, ...

Unstructured: Customer's feedback (in text), ...

# Hadoop Ecosystem

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu\_assist\_pro

# Hadoop Ecosystem – Core Components

## ❑ Hadoop Distributed File System (HDFS)

Assignment Project Exam Help

## ❑ MapReduce

<https://eduassistpro.github.io/>

Add WeChat edu\_assist\_pro



# Hadoop Distributed File System (HDFS)

- ❑ Hadoop stores files across networks using Hadoop Distributed File System (HDFS)

Assignment Project Exam Help

- ❑ Hence, Hadoop is not a **sical database**, it is a **distributed file system** (<https://eduassistpro.github.io/> s and tools in its ecosystem)

Add WeChat edu\_assist\_pro

- ❑ **Networks** can be very large, **10,000s of computers**
- ❑ HDFS is a low-level **distributed file processing system** (can be used directly for data storage)



# Hadoop Distributed File System (HDFS)

**HDFS/Hadoop** approach based on several **key assumptions**:

- ❑ **High volume:** Default physical **block sizes is 64 MB**, hence much fewer blocks per file (files are assumed to be very large)
- ❑ **Write-once, read-many:** Multiple writers can write to the same file, which improves data throughput
- ❑ **Streaming access:** Hadoop is optimized for sequential access of entire files as a continuous stream of data
- ❑ **Fault tolerance:** HDFS is designed to **replicate data** across many different devices so that when one fails, data is still available from another device (default **replication factor of three**)

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu\_assist\_pro

# Hadoop Distributed File System (HDFS)

- ❑ HDFS uses several types of nodes (computers): (see figure next slide)
  - ❑ **Data node** stores the actual file data
  - ❑ **Name node** contains file system metadata
  - ❑ **Client node** makes requests to the file system as needed to support user applications
- ❑ **Same computer** can fulfil <https://eduassistpro.github.io/nctions>
- ❑ **Data node** communicates with **name node** by sending **block** reports (list of blocks, every 6 hours) and **heartbeats** (every 3 seconds)
  - ❑ If heartbeat stops, data blocks of that node are replicated elsewhere

# Hadoop Distributed File System (HDFS)

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu\_assist\_pro

# How Does HDFS Work? [Writing]

1. The **client node** needs to create a new file, and communicates with the **name node**.
2. The **name node**
  - adds the new file name to the metadata,
  - determines a new (first) block
  - determines a list of on which data nodes to replicate,
  - and passes that information back to the client.
3. The **client node**
  - contacts the first data node specified by the name node
  - sends the data node the list of replicating data nodes.
4. First **data node** contacts the second data node in the list for replication while receiving it from the client node.
5. The **client node** gets further block numbers from the name node ... until file is written.

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu\_assist\_pro  
writing;

# MapReduce



- ☐ Implementation **complements HDFS structure**.
- ☐ **Open-source** application programming interface (API).
- ☐ **Framework** used to process large data sets across clusters.
- ☐ **“Divide and conquer”** strategy performed at node level into smaller subtasks, aggregated to final result.
- ☐ Based on **batch processing** runs tasks in parallel, ending with no user interaction.
- ☐ **YARN** (Yet Another Resource Negotiator), or MapReduce 2, can do
  - ☐ Batch processing
  - ☐ **Stream processing** (for data that comes in/out continuously)
  - ☐ **Graph processing** (for social networks)

# MapReduce

- ❑ **Map function** takes a collection of data and **sorts and filters it** into a set of key-value pairs.
  - **Mapper** program performs the map function
- ❑ **Reduce function** summarizes the map function to produce a single result.
  - **Reducer** program performs the reduce function
- ❑ **Map and reduce functions** are written as **Java** programs.
- ❑ Instead of central program retrieving the data for processing in a central location, **copies of the program are “pushed” to the nodes.**
- ❑ Typically **1 mapper *per block*, 1 reducer *per node*.**

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu\_assist\_pro

# MapReduce

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu\_assist\_pro

# MapReduce

- ❑ **Job tracker** or central control program to accept, distribute, monitor and report on jobs in a Hadoop environment
  - Typically on **name node**.
- ❑ **Task tracker** is a process responsible for reducing tasks on a node
  - Typically on **data node**.

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu\_assist\_pro



# How Does MapReduce Work?

## [Reading/Analyzing]

1. A **client node** (client application) submits a **MapReduce job** to the job tracker.
2. The **job tracker** (on server that is also the **name node**):
  - communicates with name node;
  - determines which task trackers are busy;
  - send portions of work to task trackers.
3. The **task tracker** (on server that is also a **data node**):
  - runs **map and reduce functions** (in virtual machine);
  - sends heartbeat (“still working”) and “complete” message to job tracker.
4. The **client node**
  - periodically **queries job tracker** if all task trackers are completed;
  - receives completed job.

# Hadoop Ecosystem – Data Ingestion Applications

☐ Flume

☐ Sqoop

Assignment Project Exam Help

☐ Why? **Help getting data** <https://eduassistpro.github.io/> **into Hadoop clusters.** These tools “ingest” or gather data into Hadoop.

Add WeChat edu\_assist\_pro



# Flume



- ❑ Flume is a component for **ingesting data in Hadoop**.
- ❑ Primarily for harvesting large sets of data such as **clickstream data/server logs**.
- ❑ Simple query processing **some transformation**.
- ❑ Can move data into **HDFS or HBase**.

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu\_assist\_pro

# Sqoop



- ❑ “SQL-to-Hadoop.”
- ❑ **Sqoop** is a tool for **converting data back and forth** between **relational databases** and **HDFS** (both directions).
- ❑ Works with Oracle, MySQL
- ❑ Example of Hadoop-to-SQL: MapReduce rted back into a traditional (relational) data warehouse.

<https://eduassistpro.github.io/>  
Add WeChat edu\_assist\_pro

# Hadoop Ecosystem – MapReduce Simplification Applications

❑ Hive

❑ Pig

Assignment Project Exam Help

<https://eduassistpro.github.io/>

❑ Why? **They help creating MapReduce**

- Creating MapReduce jobs requires significant skills.
- As the mapper and reducer programs become more complex, the skill requirements increase and the time to produce the programs becomes significant.

# Hive

- ❑ **Hive is a data warehousing system** that sits on top of HDFS.
- ❑ Supports its own **SQL-like language: HiveQL** (declarative / non-procedural)
- ❑ **Summarizes queries, analyzes data**

Assignment Project Exam Help

<https://eduassistpro.github.io/>

*This is the component that you use in terms of how to actually work with the data.*

Add WeChat edu\_assist\_pro



# Pig



- ❑ Hadoop platform to **write MapReduce programs**.

- ❑ Has its own **high-level scripting/programming language: Pig Latin** (procedural).

Assignment Project Exam Help

<https://eduassistpro.github.io/>

- ❑ **Pig** compiles Pig Latin scripts **into MapR** **for executing in Hadoop**.

Add WeChat edu\_assist\_pro

# Hadoop Ecosystem – Direct Query Applications

❑ HBase

❑ Impala

Assignment Project Exam Help

❑ Why? **To provide fast access to HDFS** (without going through the MapReduce processing layer)

<https://eduassistpro.github.io/>  
Add WeChat edu\_assist\_pro





# HBase



- ❑ HBase is a **NoSQL** database
- ❑ **Column-oriented** Assignment Project Exam Help
- ❑ Designed to **sit on top of** <https://eduassistpro.github.io/>
- ❑ **Quickly** processes **smaller subsets** of t Add WeChat edu\_assist\_pro
- ❑ **No SQL** support, instead uses **Java**

# Impala

- ❑ First **SQL-on-Hadoop** application
- ❑ Produced by **Cloudera**
- ❑ **SQL queries** directly against **HDFS**
- ❑ Makes heavy use of **in-memory caching**



Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu\_assist\_pro

# NoSQL Database Types

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu\_assist\_pro



# NoSQL

- ❑ **Non-relational database technologies** developed to address **Big Data challenges**
- ❑ **NoSQL** = “not modelled using relational model” (“non-SQL” / “not-only SQL”)
- ❑ **Category** emerged from organizations such as **Google, Amazon and Facebook** that **faced problems of their data sets reacting to growth**
- ❑ **Much larger data** volumes can be handled
- ❑ **Flexible structure** and often **faster**
- ❑ No standardized query language – no SQL! (maybe)
- ❑ Less adopted than RDBMS:
  - Was at peak in 2015-2016
  - Survey 2016, **16%** of companies use **NoSQL** databases and **79%** of companies use **relational databases**
- ❑ *NoSQL seems to be in decline nowadays ???!*

# NoSQL

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu\_assist\_pro

# NoSQL

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu\_assist\_pro

# NoSQL – Key-Value Database

- Store data as a collection of **key-value pairs** (keys ~ primary keys, there are no foreign keys)
- Key-value pairs are organized into logical groupings, **buckets** (<https://eduassistpro.github.io/> ~ tables)
- Key values must be unique (only) within a bucket.
- **Queries** are based on **buckets and keys** (not values)
- **get, store and delete** operations

Assignment Project Exam Help

Add WeChat edu\_assist\_pro

# NoSQL – Document Databases

- ❑ **Document databases** store data in **key-value pairs** in which the value components are **tag-encoded documents**.

Assignment Project Exam Help

- ❑ Document can be encoded in **JSON** or **BSON** (Binary JSON).
- ❑ Have tags, but still **schema-less** (not schemas, documents may have different tags).
- ❑ Documents are grouped into logical groups called **collections** (buckets).
- ❑ Tags can be queried (e.g., where balance = 0).

<https://eduassistpro.github.io/>

Add WeChat edu\_assist\_pro



# NoSQL – Column-Centric Databases

- ❑ **Column-centric (columnar) databases** focuses on storing data in columns, not rows, but still relational logic.
- ❑ **Column-centric storage:** Data in blocks which hold data from **column across many rows**
- ❑ **Row-centric storage:** Data stored in blocks which hold data from **all columns of a given set of rows**

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu\_assist\_pro

# NoSQL – Column-Centric Databases

- ❑ **Column-oriented (column family) databases** in NoSQL:
  - Organizes data in key-value pairs.
  - Keys are mapped to columns in the value component.
  - The columns vary by row.
- ❑ **Key-value pair:** name of the column and its value. Example: "cus\_Iname: Ramas".  
(~**cell** in relational model)
- ❑ **Super column:** group of columns that are logically grouped together. (positive attribute)
- ❑ **Rows keys:** created to identify objects (~**entity instances**) in the environment
- ❑ **Column family:** All of the columns (or super columns) that describe objects are grouped (~**table**)

## Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu\_assist\_pro

# NoSQL – Graph Databases

- ❑ Suitable for **relationship-rich data**

Assignment Project Exam Help

- ❑ A collection of **nodes and edges**

<https://eduassistpro.github.io/>

- ❑ **Properties** are the **attributes of a node or edge** of interest to a user

Add WeChat edu\_assist\_pro

- ❑ **Traversal** is a **query** in a graph databases

# Applications of NoSQL

- ❑ Twitter app generating 7 Tbs+ of daily tweets and displaying it back.
- ❑ Property details in a real-estate website, redundant in nature but accessed in huge numbers.  
<https://eduassistpro.github.io/>
- ❑ Online coupon sites distributing coupons.  
[Add WeChat: edu\\_assist\\_pro](#)
- ❑ Update of railway schedules and accessed by thousands of users at peak time.
- ❑ Real time score update of baseball / cricket match.

# Big Data Strategies

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu\_assist\_pro



# What is Big Data Strategy?

A Big Data strategy defines and lays out a comprehensive vision across the enterprise and sets a foundation for the organization to employ data-related or **ilities.**





Assignment Project Exam Help

<https://eduassistpro.github.io/>

Week 9 - Big Data II

Add WeChat edu\_assist\_pro

## Pre-Class Activities

-  How Do You Create A Data Strategy?
-  How to Define a Big Data Strategy
-  How to Develop a Data Strategy (Bernard Marr)
-  Types and Examples of NoSQL Databases
-  Week 9 Pre-Class Tasks

Source: <https://www.bigdataframework.org/formulating-a-big-data-strategy/>

# Challenges of Implementing Big Data Strategy

## ❑ Technological

- Lack of managerial analytics knowledge
- Technical misunderstandings and data scientists
- Inherent challenges related to data ownership and privacy
- Technical requirements in compliance regulations (e.g., NSW Transport data deluge could lead to app deluge <https://www.itnews.com.au/news/nsw-%20transport-data-liberation-could-lead-to-app-deluge-418406>)
- Costly data management tools

(Tabesh et al. 2019)



# Challenges of Implementing Big Data Strategy

## □ Cultural

- Extensive reliance on intuitive or experiential decision-making approaches
- Dominance of managerial process
- Lack of a shared understanding of its goals

Add WeChat edu\_assist\_pro  
(Tabesh et al. 2019)

### Reference:

Tabesh, P., Mousavidin, E. and Hasani, S., 2019. Implementing big data strategies: A managerial perspective. *Business Horizons*, 62(3), pp.347-358. <https://doi.org/10.1016/j.bushor.2019.02.001>

# Implementing Big Data Strategy

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu\_assist\_pro

Source: <https://www.bigdataframework.org/formulating-a-big-data-strategy/>

# Questions

Assignment Project Exam Help

<https://eduassistpro.github.io/>

Add WeChat edu\_assist\_pro

Source: stacker.com