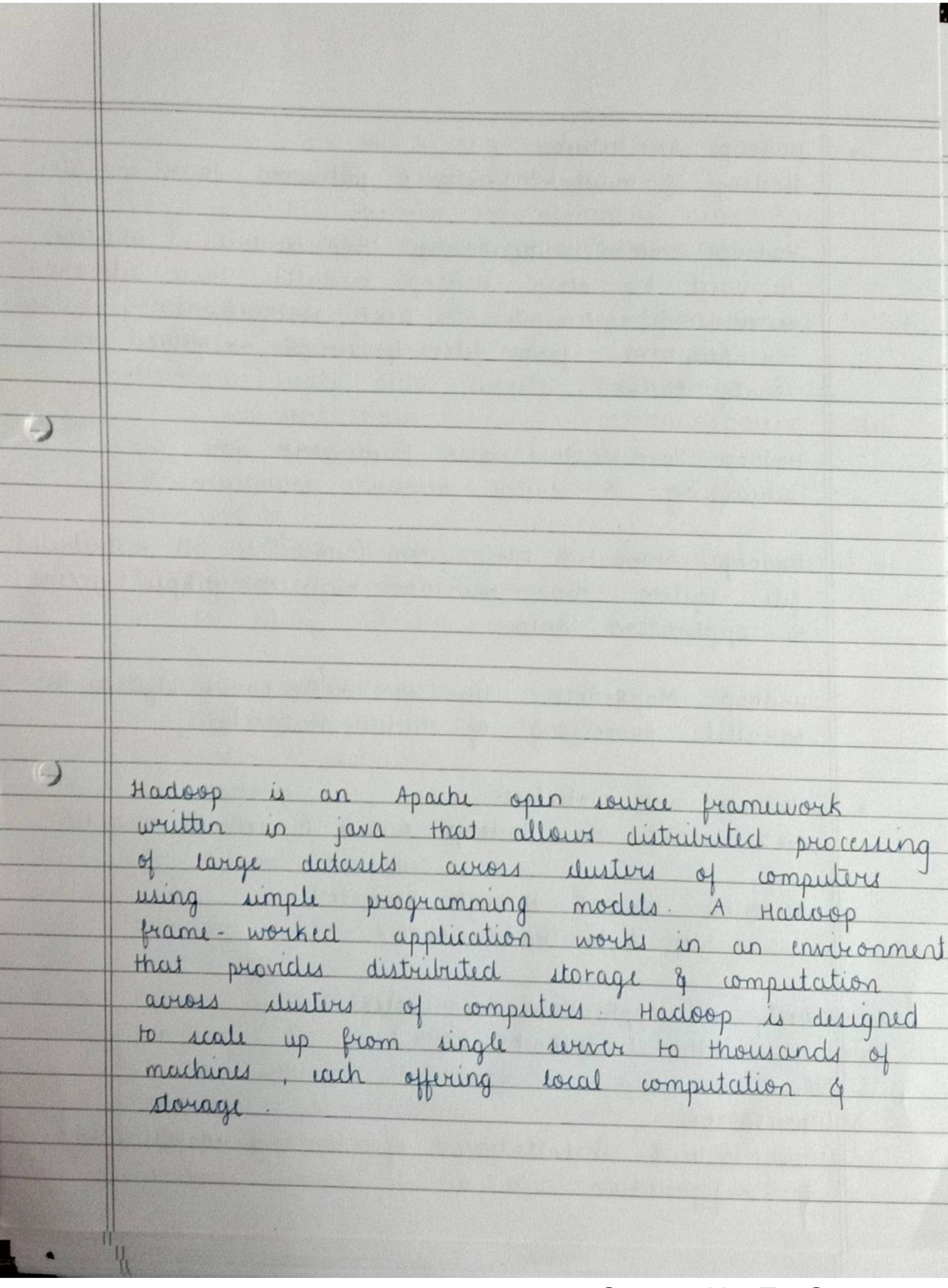
DS & BDA Lab Group A - Experiment: 1 TITLE: Single node / Multiple node Hadoop Installation:
OBJECTIVE: 1. Jo learn the concepts of Hadoop & Hadoop frame- work for Big Data. 2. Jo install & configure Hadoop.
SOFTWARE REQUIREMENTS: 1. Ubuntu 14.04 / 14.10 2. Java
Introduction: Hadoop is an open-source framework that allows to
store & process big data in a distributed environment across clusters of computers using simple programming models It is designed to scale up from single servers to thousands of machines, each offering local computation & storage. But to the advent of new technologies, devices & compounication like social networking sites, the amount of data produced by mankind is growing rapidly every year. The amount of data produced by mankind is growing rapidly every year. The amount of data produced by us from the beginning of time till 2003 was 5 billion gigabytes. The same amount was veated in every two days in 2011.

This rate is itill growing enormously. Though all this information produced is meaningful & can be wieflet when processed, it is being neglected. Big Data: Big Data means really a big data, it is a collection of large datasets that cannot be processed using traditional computing techniques. Big data is not merely a data, rather it has become a complete subject, which involves various tools, techniques & frameworks. Big data involves the data produced by different devices & applications. given below are some of the fields that come under the unbrella of Big Data. Hadoop: Doug cutting, Mike ajarella & team took the solution provided by google & started an open source Project called HADOOP in 2005 & Doug named it after his son's toy dephant. Now Apache Hadoop is a registered trademark of the Apachi software foundation. Hadoop runs applications using the Mapkeduce algorithm, where the data is processed in parallel on different CPU nodes. In short, Hadoop framework is capable enough to develop appli-- cations capable of running on clusters of computers & they could perform complete statistical analysis for a huge amount of data.



Hadoop Architecture: Hadoop framework includes following four modules:
Hadoop common: These are gava libraries & utilities required by other Hadoop modules. These libraries provide filesystem & os evel abstractions & contains the necessary Java files & scripts required to start Hadoop.
Hadoop Yarn: This is a framework for job- schiduling & duster resource management.
Hadoop Distributed File System (HDFS TM): A distributed file system that provides high throughput accers to application data.
Hadoop MapReduce: This is YARN-based system for parallel processing of large datasets.
Hadoop web Interfaces: Let's start the Hadoop again & see its web UI:
Accessing HADOOP through browser: http://localhost:50070/
Verify all applications for clusters: http://localhost:8088/
CONCLUSION: We studied installation of Hadoop installation & configuration.

Group B - Experiment 4 TITLE: Perform the operations using Python on the Facebook metrics data sets
AIM/: Perform the following operations using Python- PROBLEM on the Jacobook metrics data sets: STATEMENT:(a) (reate data subsets (b) Merge Data (c) sort Data (d) Iransporing Data (e) shape & rushape Data
OBJECTIVE: 1. Jo undvuland & apply the Analytical concept of Big data using Python. 2. To study detailed concept Python.
SOFTWARE REQUIREMENTS: 1. Ubuntu 14.04/14.10 2. GNU (Compiler 3. Hadoop 4. Java 5. Python Platform.
Python is an easy to learn, powerful programming language It has efficient high - level data structures of a simple but effective approach to object-oriented programming Python's eligant syrdax & dynamic typing, together with its interpreted notive, make it an ideal language for

	areas on most platforms The Python interpreter & the extensive it and area library are freely available in source or binary form for all major platforms from the Python web site, https://www.python.org/, & may be freely distributed. The same site also contains distributions of & pointers to many free third-party Python modules, programs & tools, & additional documentation.
	Features of Python: Python is a dynamic, high livel, free open source & interpreted programming language. i) Early to code ii) object - oriented language. iii) Tree & open source iv) GUI Programming support y High livel language vi) Extensible feature vii) Python is portable.
*	Data Reshaping: Data Reshaping is about changing the way data is organised into rows & columns. Most of the time data processing in python is done by taking the input data as data frame. It is easy to extract data from rows & columns of data frame but there are situations when we need the data frame in a format that is different from format in which we received it. Python has many functions to split, merge &

	g change the rows to columns & via versa is a data frame.
1)	Joining columns & Rows m a Data Frame: We can join multiple vectors using chind() function. We can murge using rhind() function.
*	Merging Data frames: We can merge dataframes by using merge() function. We can add columns using merge & rows using rhind functions.
2)	Subsets: subsetting Vectors, Matrices & Data Frames. Return subsets of rectors, matrices or data frames which meet conditions.
5)	Arguments: x: object to be subsetted: subset: logical expression indicating elements or nows to keep; missing values are taken as false. select: expression, indicating columns to select from a data frame. drop: passed on to indexing operatore.
3)	Melting & casting: The functions used to do this are called melt() & cast (). Steps:- Melt the data, cast the Molten Data, Transpose using t() function.

4)	Sorting Data:
	To sort a dataframe use order() function.
	By default, sording is ASCENDING Prepand
	the sorting variable by a minus sign to indicate DESCENDING order.
	DESCENDING OGGG
	CONCLUSION: They was hard learnet hours to neck
	conclusion: Thus we have learnt how to perform the different reshape operations using python.
	J Pyrram