

Big Data Analytics

COMP336 / COMP529

Dom Richards

University of Liverpool

November 11, 2024

Table of Contents

- 1 Introduction
- 2 Course Overview
- 3 Distributed File Systems and MapReduce
- 4 Spark
- 5 Classification
- 6 Machine Learning for Text Analysis

Introduction

1

Introduction

- Characteristics of Big Datasets
- Dataset Representations

Attendance: 500240

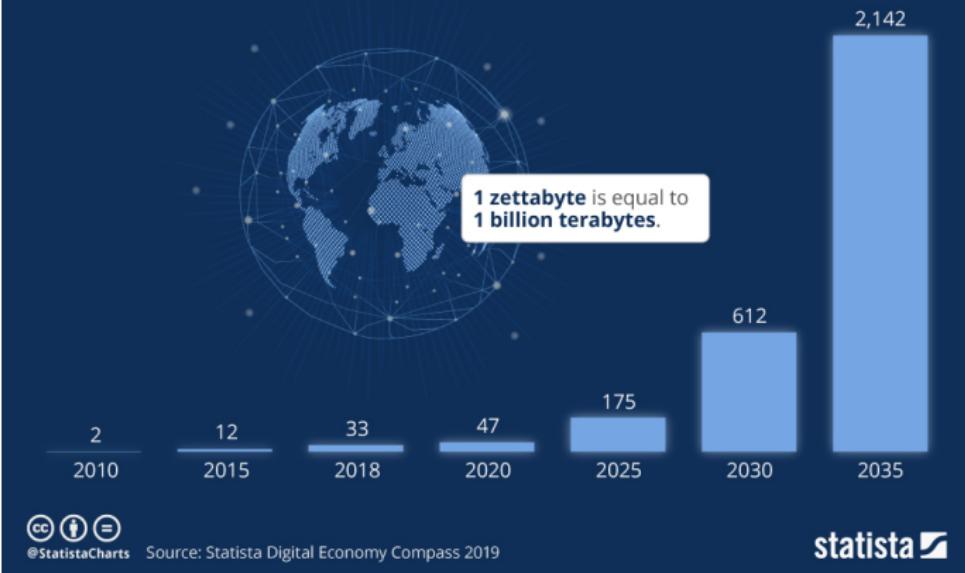
Introduction

- Welcome to Big Data Analytics!
- Over the next 12 weeks, we'll explore the field of big data
- It's concerned with the storage and analysis of datasets that are too large to fit on a computer, and hence typically require a *data center* to store and process
- 'Big' datasets (also referred to as big data) are often very valuable, and they crops up everywhere: written text; social media posts; financial transaction records; industrial sensors; scientific research data
- By the end of this course, you'll have a solid understanding of big data technologies and how to apply them to real-world problems
- We'll introduce a range of techniques to extract value from big datasets including: supervised and unsupervised machine learning; statistical analysis; natural language processing; network science; signal processing

How Common is Big Data?

Global Data Creation is About to Explode

Actual and forecast amount of data created worldwide 2010-2035 (in zettabytes)



Source [Statista report on global data creation forecasts](#).

Example Datasets

- Search Engine Data: googlebot
- Social Media Data: Posts, comments, and reactions on platforms like Facebook and Twitter
- E-commerce Data: Online shopping behaviours, product views, and purchases
- Healthcare Data: Patient records, genetic sequences, and medical device data
- IoT Data: Sensor data from connected devices like smart homes and wearables

1

Introduction

- Characteristics of Big Datasets
- Dataset Representations

Attendance: 500240

Understanding Big Datasets

- In the next two sections, we'll explore the characteristics of big datasets and the challenges they present
- In this section, we'll discuss the 6 Vs of big data: volume; velocity; variety; veracity; value; and variability
- In the next section, we'll take a look at how big data is represented for the purpose of analysis

Characteristics of Big Data

- Big data comes in many shapes and sizes, but it can be characterised in terms of “6 Vs”:
 - Volume: the amount of data in a given dataset
 - Velocity: Data can be generated at a range of possible speeds
 - Variety: Data comes in different formats, including structured, semi-structured, and unstructured
 - Veracity: Data quality and accuracy can vary
 - Value: which derives from the ability to transform big data into actionable insights
 - Variability: Data can be inconsistent or change rapidly

Volume

- Volume refers to the quantity of data
- Ranges from GBs upwards

Velocity

- Velocity refers to the speed at which data is being generated, processed, and analysed
- For some big data applications, data arrives continuously from multiple sources in real time
- Example: Credit card transactions need to be validated quickly to prevent fraud
- Rapid decision-making is crucial in fields like:
 - Stock trading platforms, where delays can result in losses
 - Autonomous vehicles, which rely on sensor data being processed in real-time

Variety

- Variety refers to the wide range of data types and formats
- Traditional systems only handled structured data, typically found in relational databases (e.g., SQL)
- Today, data comes in a variety of formats:
 - Structured data: Databases, spreadsheets, CSV, Parquet
 - Unstructured data: Emails, videos, social media posts, images

Veracity

- Veracity refers to the accuracy, quality, and trustworthiness of data
- Big data often includes noisy, incomplete, or inconsistent data
- Example challenges:
 - Social media data may include misleading or biased information
 - Sensor data can contain errors due to faulty devices
- Techniques to improve veracity:
 - Data cleaning: Removing inaccuracies or noise
 - Data validation: Ensuring data is correct before analysis

Value

- Value refers to the insights and actionable information extracted from big data
- Big data is often valuable if it can be transformed into meaningful insights
- In this course, we'll see a range of techniques to extract value from big data, including: machine learning; statistical analysis; natural language processing; network science; signal processing.

Variability

- Variability refers to the changing nature of data and the inconsistency in its flow
- Data can be generated at unpredictable rates and can vary in its meaning and relevance
- E.g. regime switching:
 - Social media sentiment analysis can shift rapidly based on trending topics or events
 - IoT sensor data might fluctuate based on environmental changes or technical issues
- Handling variability requires adaptive algorithms and real-time monitoring to respond quickly to data changes

1

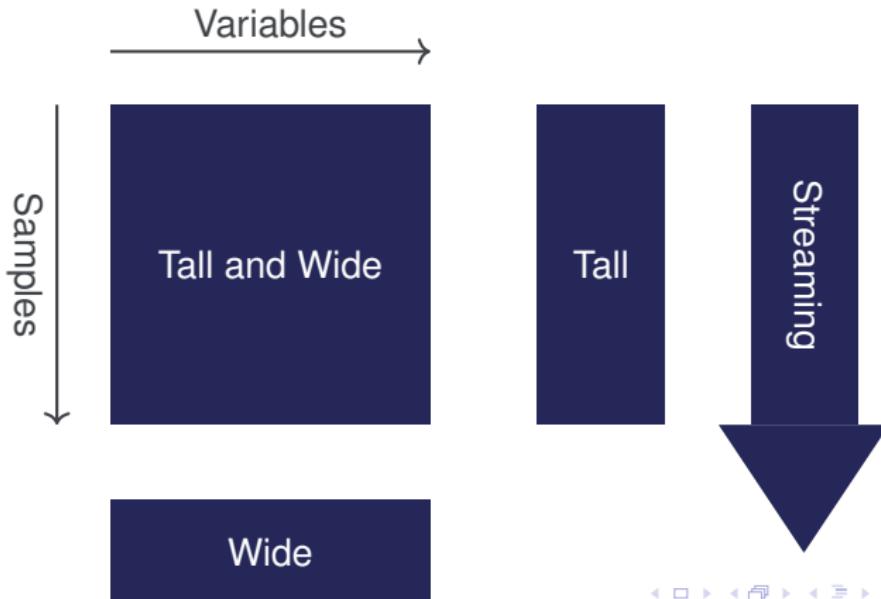
Introduction

- Characteristics of Big Datasets
- **Dataset Representations**

Attendance: 500240

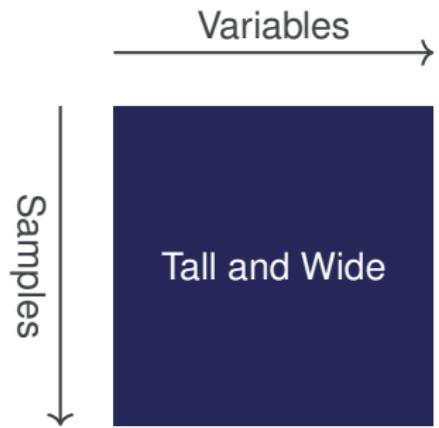
Dataset Types

Big datasets come in various shapes and sizes. We tend to think about matrices of variables and samples, even when the data isn't originally presented to us like this (e.g. text documents or images)



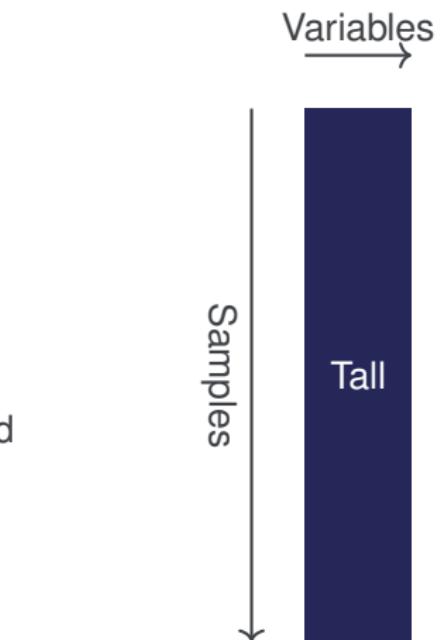
Tall and Wide Datasets

- These datasets can be very large:
 - Millions to billions of variables and samples
- Typical example is networks:
 - Social networks, contagion networks, etc.



Tall Datasets

- Hundreds to thousands of variables
- Thousands to millions of samples
- Often suitable for machine learning
- Examples:
 - NLP (Natural Language Processing)
 - Image databases
 - Scientific computing (e.g., solutions of fluid dynamics equations)
 - Reinforcement learning



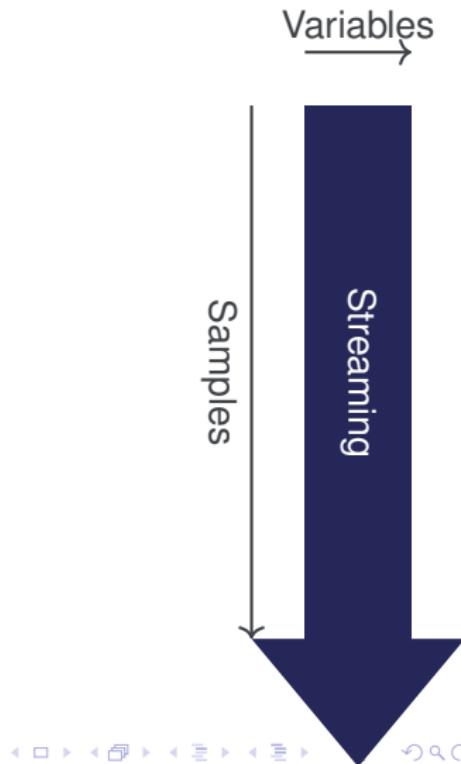
Wide Datasets

- Thousands to millions of variables
- Hundreds to thousands of samples
- Variables \gg samples
- Examples:
 - Genomics datasets
 - Document modelling



Streaming Datasets

- Samples continually created in real time
- Examples:
 - Sensor data
 - Social media feeds
 - Financial transactions



Dense vs Sparse Matrices

- Data matrices can be dense or sparse
- Dense matrix:
 - Most matrix elements are non-zero
 - All elements are stored, including the zeroes
 - Memory usage is proportional to the total number of elements (*variables* \times *samples*)
- Sparse matrix:
 - Majority of elements are zero
 - Only non-zero elements are stored, often as a coordinate list, which stores triplets (row, column, value)
 - Memory usage is proportional to the number of non-zero elements

Course Overview

2

Course Overview

- Big Data Middleware
- Network Science
- Unsupervised Learning
- Supervised Learning
- Natural Language Processing
- Bayesian Statistics and Filters

Attendance: 500240

2

Course Overview

● Big Data Middleware

- Network Science
- Unsupervised Learning
- Supervised Learning
- Natural Language Processing
- Bayesian Statistics and Filters

Attendance: 500240

Working with Big Data

- In weeks 2, 3 and 11, you'll learn about big data middleware
- Datasets in big data are typically too large to be processed on a single computer, so software has been developed to process such datasets across multiple computers in a data centre
- This software is called *big data middleware*
- It is installed onto a cluster of compute nodes in a data centre, and allows users to interact with the entire cluster as if it were a single computer
 - Storage is distributed across the cluster, and the user is presented with a single file system
 - Processing is distributed across the cluster, allowing the user to operate on all of the data by writing a single program which is executed on all the cluster nodes in parallel
- You'll learn about the two most popular open source big data middleware frameworks: Hadoop and Spark

Hadoop

- Next week we'll see Hadoop.
- Hadoop was the first major open-source big data middleware framework
- It was based on Google's MapReduce and Google File System papers (which are required reading for next week)
- The core components of Hadoop:
 - HDFS (Hadoop Distributed File System): A scalable and fault-tolerant file storage system
 - MapReduce: A programming model for processing large datasets in parallel. Performs fairly simple operations (map and reduce) at huge scales

Storing Data with HDFS

- The Hadoop Distributed File System (HDFS) provides scalable, fault-tolerant data storage
- HDFS splits large files into blocks and distributes them across multiple machines for parallel storage and retrieval

Spark

- In weeks 3 and 11, we'll look at Spark
- Spark was developed after Hadoop to address some of its perceived limitations
- It's often faster than Hadoop (due to a more sophisticated use of RAM on cluster nodes), and can perform a broader range of tasks
- You'll learn about Spark in week 3
- In week 11, you'll see how Spark can perform advanced data analysis, including machine learning and graph processing:
 - Spark's MLLib library supports a wide range of machine learning algorithms
 - Spark's GraphX library enables complex graph computations

2

Course Overview

- Big Data Middleware
- **Network Science**
- Unsupervised Learning
- Supervised Learning
- Natural Language Processing
- Bayesian Statistics and Filters

Attendance: 500240

Network Science

- In weeks 4, you'll move on to the task of extracting value from big datasets
- You'll learn about Network Science, and its use in processing network (AKA graph) based big datasets
- A network consists of nodes (or vertices) connected by edges (or links)
- Examples of networks in big data include social networks, biological networks, and transportation systems
- Network Science studies the structure, dynamics, and behaviour of networks

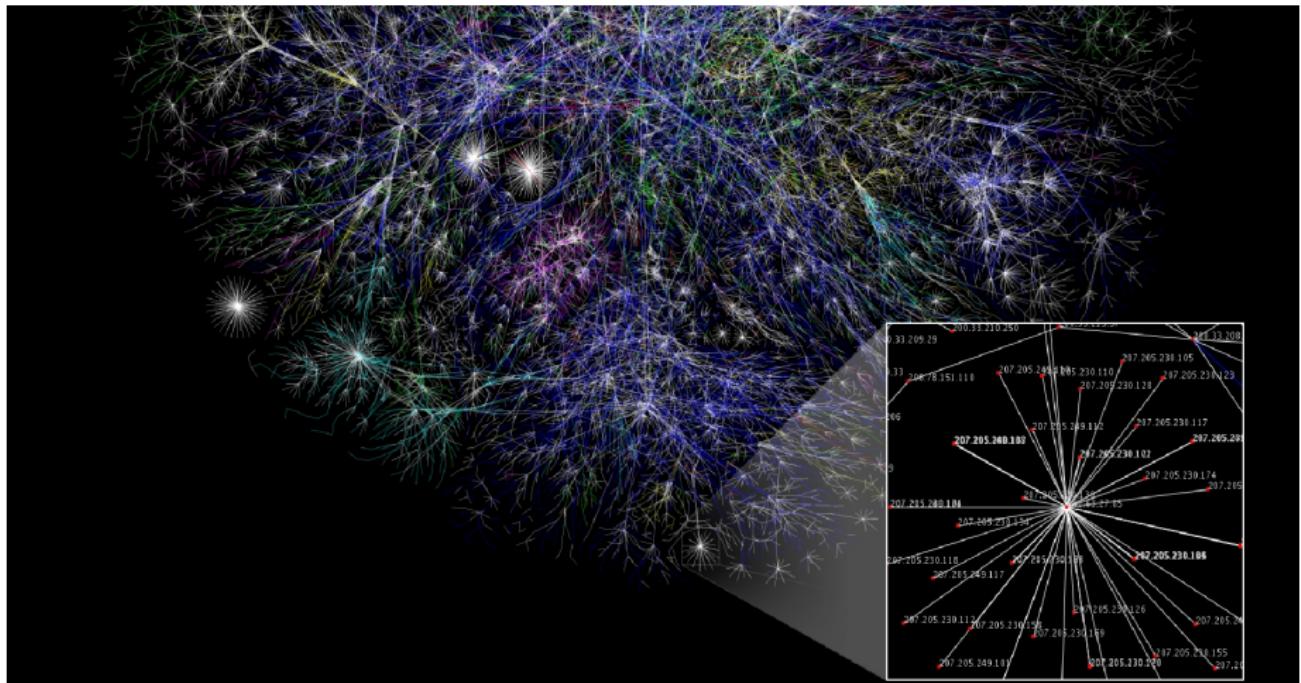
Facebook Connections



Facebook Connections Map the World

- The image shows a map of Facebook connections as of 2010
 - The lines on the map don't represent physical features like coasts or borders, but instead reflect real human relationships and social networks
 - The map also highlights areas with minimal Facebook presence c. 2010, such as China and central Africa, where few connections exist

Internet Map



Internet Map Visualisation

- This image shows a map of the Internet, visualising the global structure of interconnected networks
- Each node represents a network, and the edges represent the connections between these networks
- Such visualisations help understand the topology and connectivity of the Internet infrastructure

Network Representations

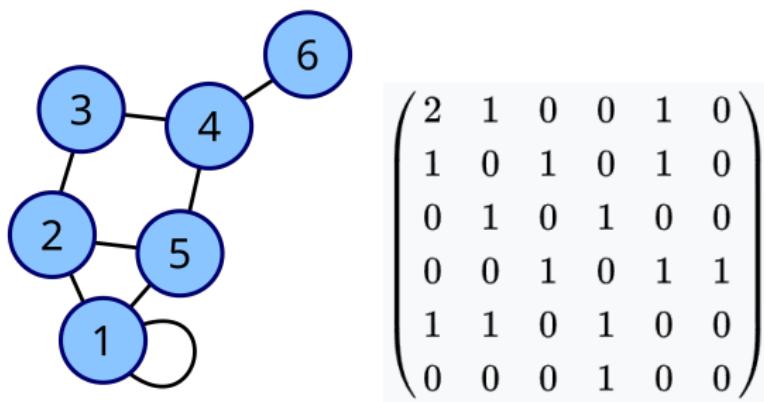
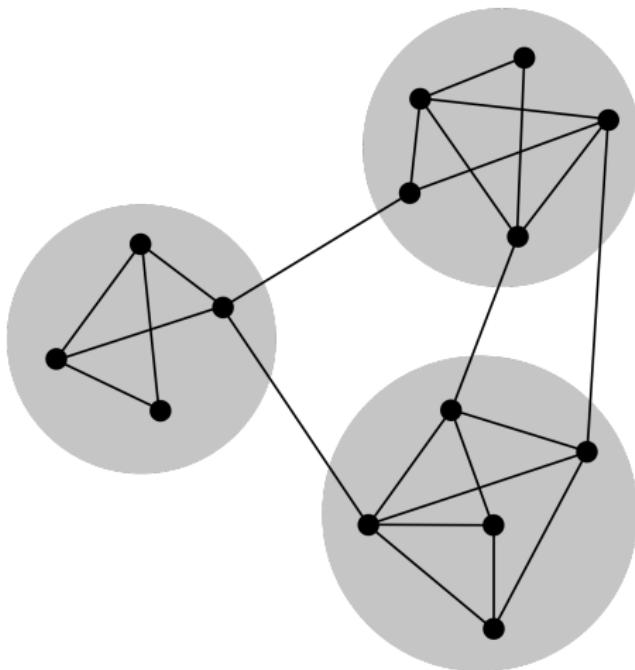


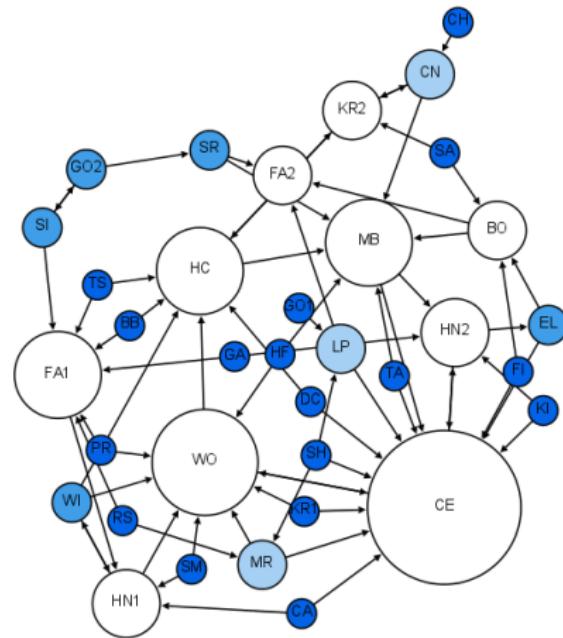
Figure: Comparison between graph and its adjacency matrix.

Network Community Structures



Preferential Attachment in Networks

The Moreno sociogram of a 1st grade class:



Scale-Free Networks

- **Scale-free networks** are a type of network characterised by a power-law degree distribution
- In these networks, most nodes have few connections, while a few nodes (called "hubs") have many connections
- They have a **Power-law distribution**, as you'll see in week 4
- Examples of scale-free networks include the internet, social networks, and biological systems like protein interactions

2

Course Overview

- Big Data Middleware
- Network Science
- **Unsupervised Learning**
- Supervised Learning
- Natural Language Processing
- Bayesian Statistics and Filters

Attendance: 500240

Unsupervised Learning: PCA and Clustering

- In weeks 5 and 6, you'll learn about unsupervised learning
- Unsupervised learning deals with unlabelled data, where the goal is to find hidden patterns or structure in the data
- You'll see two of the most common techniques in unsupervised learning:
 - **Principal Component Analysis (PCA)** (week 5):
 - A dimensionality reduction technique that projects data onto the principal components—directions of maximum variance
 - PCA helps simplify datasets by reducing features while retaining most of the important information
 - **Clustering** (week 6):
 - Clustering algorithms, such as K-means and hierarchical clustering, group similar data points into clusters based on their features
 - The goal is to partition data into distinct groups where points within a group are more similar to each other than to those in other groups

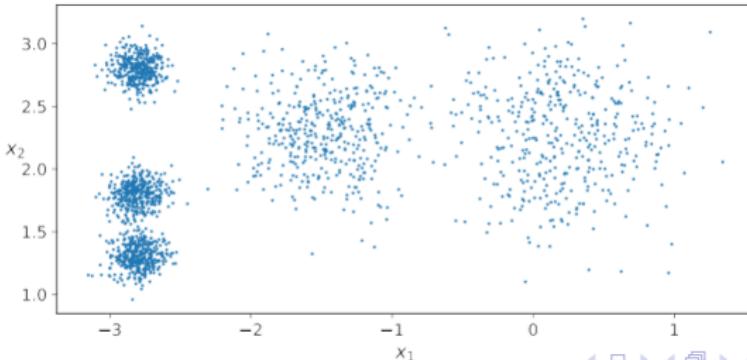
PCA for Exploratory Data Analysis

- We can use PCA to reduce the dimensionality of data all the way down to 2 or 3 dimensions
- At that point, the data can be visualised
- E.g., the plot shows the first two principal components for the handwritten digits dataset
- This often yields insights that can be exploited with further application of machine learning



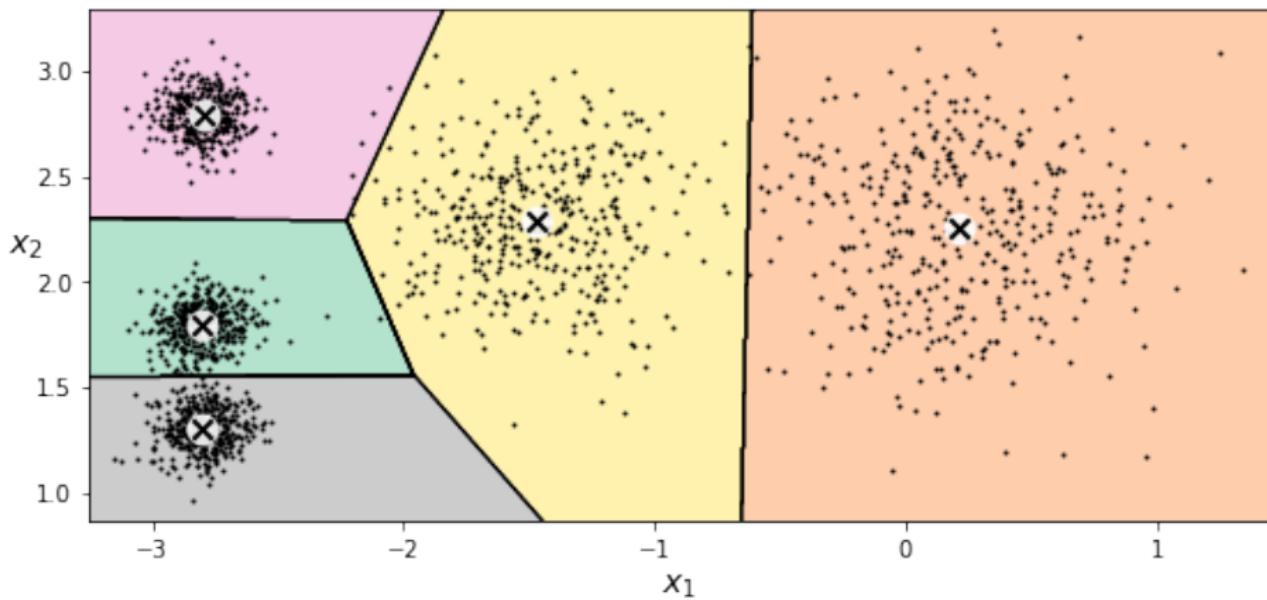
Clustering

- Clustering is an unsupervised learning technique that groups similar data points into distinct clusters
- Clustering allows us to group data points based on their similarity or distance from each other
- It's useful for exploratory data analysis, pattern recognition, and anomaly detection
- Popular algorithms include K-Means, DBSCAN, and Hierarchical Clustering, all of which you'll see in week 6

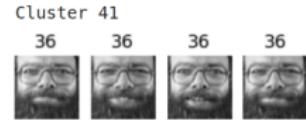
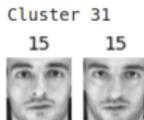


Clustering with K-Means

Here's the result of applying the K-Means clustering algorithm to the above dataset:



Clustering Faces



Clustering-Based Image Compression

Original image



10 colors



8 colors



6 colors



4 colors



2 colors



2

Course Overview

- Big Data Middleware
- Network Science
- Unsupervised Learning
- Supervised Learning**
- Natural Language Processing
- Bayesian Statistics and Filters

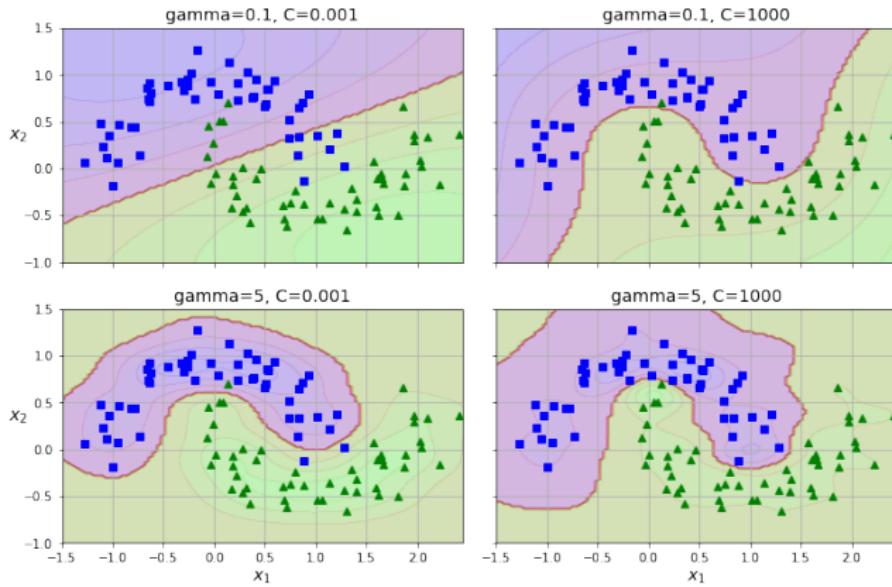
Attendance: 500240

Supervised Learning

- In week 8, we'll look at supervised learning
- Supervised learning is a type of machine learning where the model is trained on labelled data
- The model learns to map input data to output data based on the input-output pairs provided during training
- We'll look at several supervised learning algorithms, including k-Nearest Neighbours and Support Vector Machines

Half Moon Classification

The plot illustrates the use of Support Vector Machines (SVM) on a dataset



2

Course Overview

- Big Data Middleware
- Network Science
- Unsupervised Learning
- Supervised Learning
- **Natural Language Processing**
- Bayesian Statistics and Filters

Attendance: 500240

Natural Language Processing

- In week 8, we'll also look at Natural Language Processing (NLP), which enables computers to understand, interpret, and respond to human language
- We'll focus on the task of document search, where the goal is to retrieve relevant documents from a corpus, usually based on a user's query
- We'll look at “bag of words” style document representations, culminating in an approach called Term Frequency-Inverse Document Frequency (TF-IDF):

2

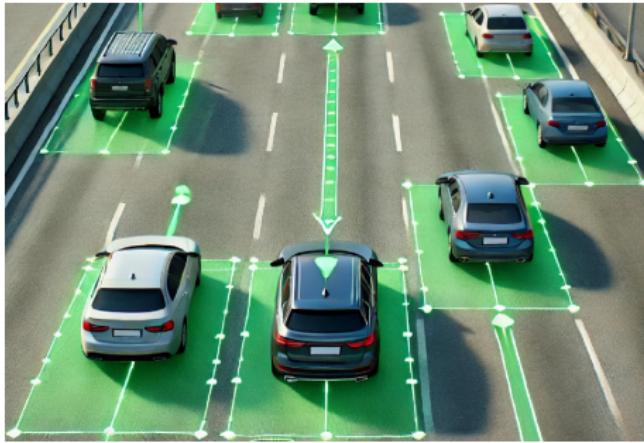
Course Overview

- Big Data Middleware
- Network Science
- Unsupervised Learning
- Supervised Learning
- Natural Language Processing
- **Bayesian Statistics and Filters**

Attendance: 500240

Bayesian Statistics and Filters

- In weeks 7, 9 and 10, you'll learn about Bayesian statistics and filtering
- This will culminate with an introduction to the Kalman filter - a very widely-used algorithm that uses a series of measurements observed over time to estimate the current state of a system and predict its future evolution
- Applications include vehicle tracking, robotics and navigation



Course Outline

- ① Introduction and Course Overview
- ② Hadoop
- ③ Spark
- ④ Network Science
- ⑤ Linear Algebra Approaches
- ⑥ Clustering
- ⑦ Probabilistic Modelling
- ⑧ Classification and Real World Applications
- ⑨ Parallel Computing, Bayesian Modelling and Real World Applications
- ⑩ Parallel Computing, Bayesian Modelling and Real World Applications
- ⑪ Advanced Spark
- ⑫ Advanced Approaches

Prep for Next Week

Please read through Google's original white papers on Google File System (GFS) and MapReduce:

- **The Google File System, 2003**
- **MapReduce: Simplified Data Processing on Large Clusters, 2004**

They're very accessible, and they laid the groundwork for Hadoop and HDFS

The Google File System

Sanjay Ghemawat, Howard Gobioff, and Shun-Tak Leung
Google

ABSTRACT

We have designed and implemented the Google File System, a scalable distributed file system for large distributed data sets. It is designed to run on commodity hardware running on inexpensive commodity hardware, and it delivers high aggregate performance to a large number of clients. While sharing many of the same goals as previous distributed file systems, such as availability and scalability, our design is informed by lessons learned from observations of our application workloads and technological environment, both current and anticipated, that reflect a marked departure from earlier distributed file systems. Our design has led to more than a dozen novel choices and explores radically different design points.

This paper describes how we successfully built our storage needs. It is widely deployed within Google as the storage platform for the generation and processing of data used by one service in size and complexity. The system has been deployed to handle large data sets. The largest cluster to date contains hundreds of terabytes of storage across thousands of disks on over a thousand hosts. This paper also discusses some of the challenges we faced in building a system that can be easily scaled by both size and complexity.

In this paper, we describe the design of GFS, and inspect its multi-level architecture.

Categories

D.4 [3 – Datal]

General Terms

Design, reliability

Keywords

Fault tolerance

[✉]The authors: {sanjay, hgoobl}@

Copyright 2003

MapReduce: Simplified Data Processing on Large Clusters

Jeffrey Dean and Sanjay Ghemawat

jeff@google.com, sanjay@google.com

Google, Inc.

Abstract

MapReduce is a programming model and an associated implementation for processing and generating large parallelizable key/value pairs. Users specify a map function that processes a key/value pair to generate a set of intermediate key/value pairs, and a reduce function that merges all intermediate values associated with the same intermediate key. Many real world tasks are expressible in this model, as shown in the paper.

MapReduce writes to the functional style automatically parallelized and executed on a large cluster of commodity machines. The run-time system takes care of the details of partitioning the input data, scheduling the program's execution across a set of machines, handling machine failures, and managing the required inter-machine

give day, etc. Most such computations are conceptually straightforward. However, the input data is usually large and the computations have to be distributed across hundreds or thousands of machines. In order to keep it simple, MapReduce provides a small set of primitives. The interface is designed to let users focus on the logic of the algorithm, not the computation, distribution, data, and handle failures compute to obscure the original simple computation with large amounts of complex code to deal with these issues.

As a reaction to this complexity, we designed a new abstraction to make it easier to write distributed programs. Specifically, we were trying to implement but failed the necessary details of parallelization, fault-tolerance, data distribution and load balancing in a library. Our abstraction is inspired by the map and reduce primitives present in Lisp and many other functional languages. We realized that

Distributed File Systems and MapReduce

3

Distributed File Systems and MapReduce

- Google MapReduce and Google File System
- Motivating Examples from Search Engine Design
- Hadoop
- The Hadoop Distributed File System (HDFS)
- Orchestration of Hadoop MapReduce

Attendance: 500240

Big Data Middleware

- Today and in weeks 3 and 11, we'll focus on middleware for big data analysis
- Many datasets (or collections of datasets) are too big to fit onto a single computer
 - We'll see a few examples today
- These datasets are stored in a distributed way across different servers in a data centre
- Big data middleware sits on top of data centre hardware. It:
 - Handles the distributed storage of data
 - Provides a computational abstraction (AKA programming model) to make the storage and subsequent processing of these datasets more tractable

Big Data Middleware

- We'll start with the original big data middleware: Google File System (GFS) and Google MapReduce, which were the first big data middleware solutions
 - Google published papers on them in 2003 and 2004 respectively (and implemented them internally prior to this)
 - We'll introduce them with the original motivating examples from Google
 - We'll see that "map reduce" is a very simple computational paradigm, but you can do a lot of useful big processing with it
 - Google unfortunately hasn't released the source code for GFS and MapReduce, so you can't run it

Big Data Middleware

- We'll then look at the first open source big data middleware solution, which was Hadoop MapReduce and the Hadoop Distributed File System (HDFS)
- All of this will lay the groundwork for introducing Apache Spark next week, which is a fast, cutting edge generalisation of the MapReduce concept

Required Reading for Today's Lecture

- The Google File System, 2003
- MapReduce: Simplified Data Processing on Large Clusters, 2004

The Google File System

Sanjay Ghemawat, Howard Gobioff, and Shun-Tak Leung
Google

ABSTRACT

We have designed and implemented the Google File System, a scalable distributed file system for large distributed data sets. It is designed to run on inexpensive commodity hardware, and it delivers high aggregate performance to a large number of clients.

While sharing many of the same goals as previous distributed file systems, such as performance, scalability, reliability, and availability, its design has been informed by lessons learned from real-world loads and technological advancements, both current and anticipated, that reflect a marked departure from earlier system design assumptions. We have reexamined traditional designs and explored radically different points in the design space.

This paper describes our design and early experiences running tens of thousands of storage nodes. It is widely deployed within Google as the storage platform for the generation and processing of data used by one service or more to support a range of applications that handle large data sets. The largest cluster to date contains hundreds of terabytes of storage across thousands of disks on over a thousand hosts.

The system is designed to run on commodity hardware and is built from inexpensive components.

In this paper, we describe the design and implementation of our system and discuss some of the interesting issues we faced.

Categories
D [4]: 3—Data

General Terms
Design, reliability

Keywords
Fault tolerance
MapReduce
The authors
[sanjay, hgs]@google.com

MapReduce: Simplified Data Processing on Large Clusters

Jeffrey Dean and Sanjay Ghemawat

jeff@google.com, sanjay@google.com

Google, Inc.

Abstract

MapReduce is a programming model and an associated implementation for processing and generating large data sets. Users specify a map function that processes a key/value pair to generate a set of intermediate key/value pairs, and a reduce function that merges all intermediate values associated with the same intermediate key. Many real world tasks are expressible in this model, as shown in the paper.

Given the simplicity of the functional style we automatically parallelized and executed on a large cluster of commodity machines. The run-time system takes care of the details of partitioning the input data, scheduling the program's execution across a set of machines, handling machine failures, and managing the required inter-machine communication.

Given the complexity of the computation, most such computations are conceptually straightforward. However, the input data is usually large and the computations have to be distributed across hundreds or thousands of machines. In order to keep in memory the intermediate data, we had to make some compromises. This involves the need to copy data, to allocate the computation, distribute the data, and handle failures compute to obscure the original simple computation with large amounts of complex code to deal with these issues.

As a reaction to this complexity, we designed a new abstraction that allows us to ignore many of the details we were trying to implement but hide the messy details of parallelization, fault-tolerance, data distribution and load balancing in a library. Our abstraction is inspired by the map and reduce primitives present in Lisp and many other functional languages. We realized that

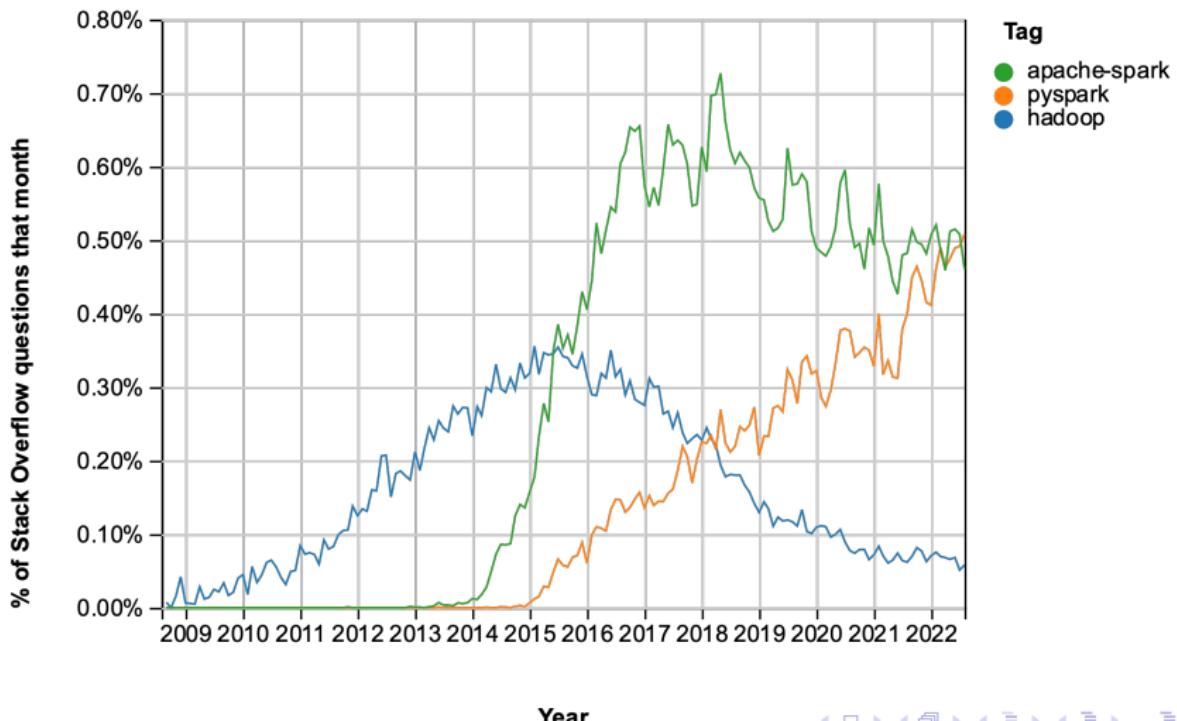
Big Data Middleware at Google

- The original example of big data middleware was Google's proprietary MapReduce and Google File Store:
 - Described by Google in 2003 and 2004 research papers, and implemented at Google some years before that
 - It provided a computational abstraction: simple distributed processing + result aggregation, as we'll see
 - It was designed for use by Google's search engine
- Google's MapReduce software is proprietary and has never been released to the public, but shortly after the 2003 and 2004 publications, Hadoop and HDFS emerged as a free open source equivalent, which use the same MapReduce computational abstraction

Hadoop and Spark

- We'll focus on Spark in week 3, and we'll see advanced Spark algorithms in week 11
- Spark is a newer alternative to Hadoop
 - It can perform the same MapReduce tasks as Hadoop, often with a speed-up of 10x to 100x
 - It also performs a more general set of computations which generalise MapReduce, whilst maintaining the same high performance

Hadoop and Spark Popularity Over Time



3

Distributed File Systems and MapReduce

- **Google MapReduce and Google File System**
- Motivating Examples from Search Engine Design
- Hadoop
- The Hadoop Distributed File System (HDFS)
- Orchestration of Hadoop MapReduce

Attendance: 500240

Google's Datasets

- Google builds its search engine by using a piece of software known as a "web crawler" to explore the entire web continually in order to find pages, which are then added to a giant index
 - As you might imagine, this generates a lot of data
- The web crawler is called Googlebot
- For more information, see:
 - [https://developers.google.com/search/docs/fundamentals/
how-search-works](https://developers.google.com/search/docs/fundamentals/how-search-works)
 - [https://developers.google.com/search/docs/
crawling-indexing/googlebot](https://developers.google.com/search/docs/crawling-indexing/googlebot)

Datasets and Data Centre Node Specs

- From Google's MapReduce paper, as of 2004:
 - Google had multiple data centres
 - Each data centre had hundreds or thousands of servers
 - Hard disk capacity per server: 320GB
 - Memory per server: 4GB
 - Network bandwidth per server: 100-200Gbps
 - Average file size: 113GB
- Data was largely generated by the Googlebot web crawler, as we'll see
- Note that individual files were usually small enough to fit onto a single server (113GB vs 320GB).
 - However, the entire file store was far too large to fit onto a single server

Datasets and Data Centre Node Specs

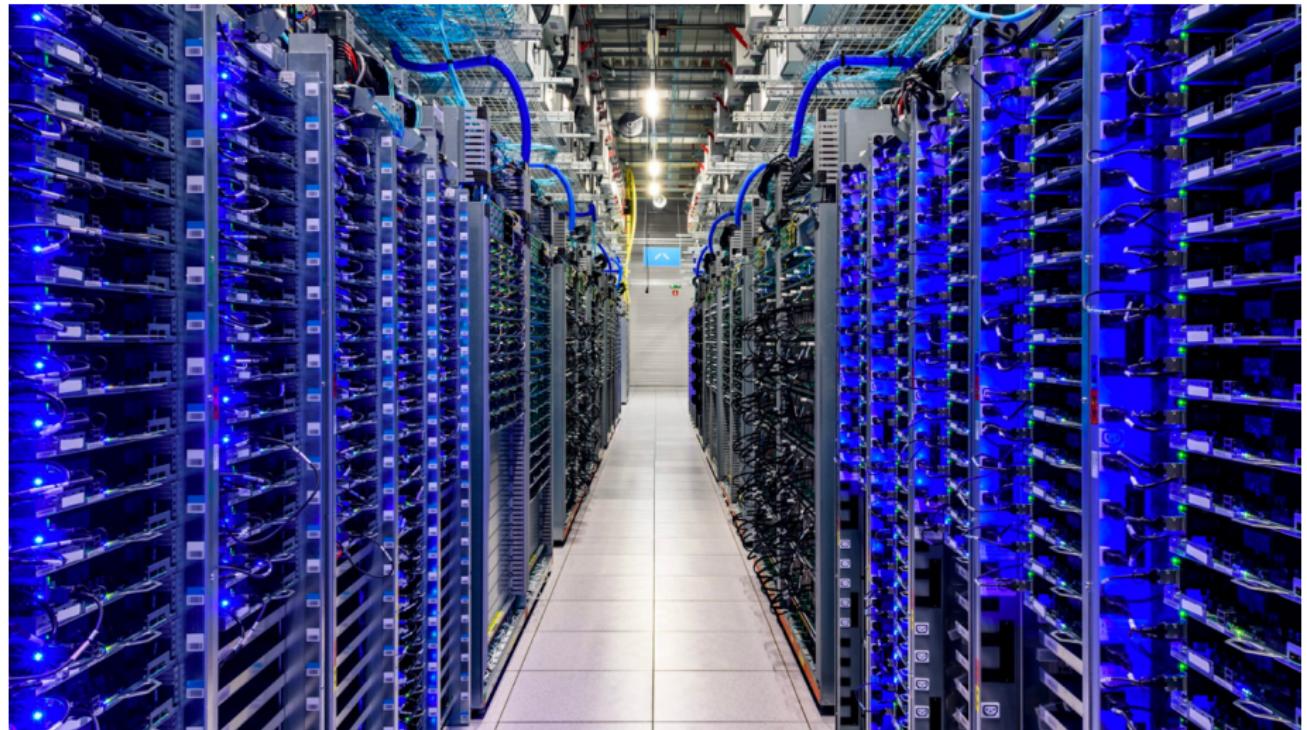
| | |
|---------------------------------------|-------------|
| Number of jobs | 29,423 |
| Average job completion time | 634 secs |
| Machine days used | 79,186 days |
| Input data read | 3,288 TB |
| Intermediate data produced | 758 TB |
| Output data written | 193 TB |
| Average worker machines per job | 157 |
| Average worker deaths per job | 1.2 |
| Average map tasks per job | 3,351 |
| Average reduce tasks per job | 55 |
| Unique <i>map</i> implementations | 395 |
| Unique <i>reduce</i> implementations | 269 |
| Unique <i>map/reduce</i> combinations | 426 |

Table 1: MapReduce jobs run in August 2004

AWS Data Centre External View



Google Data Centre Inside View



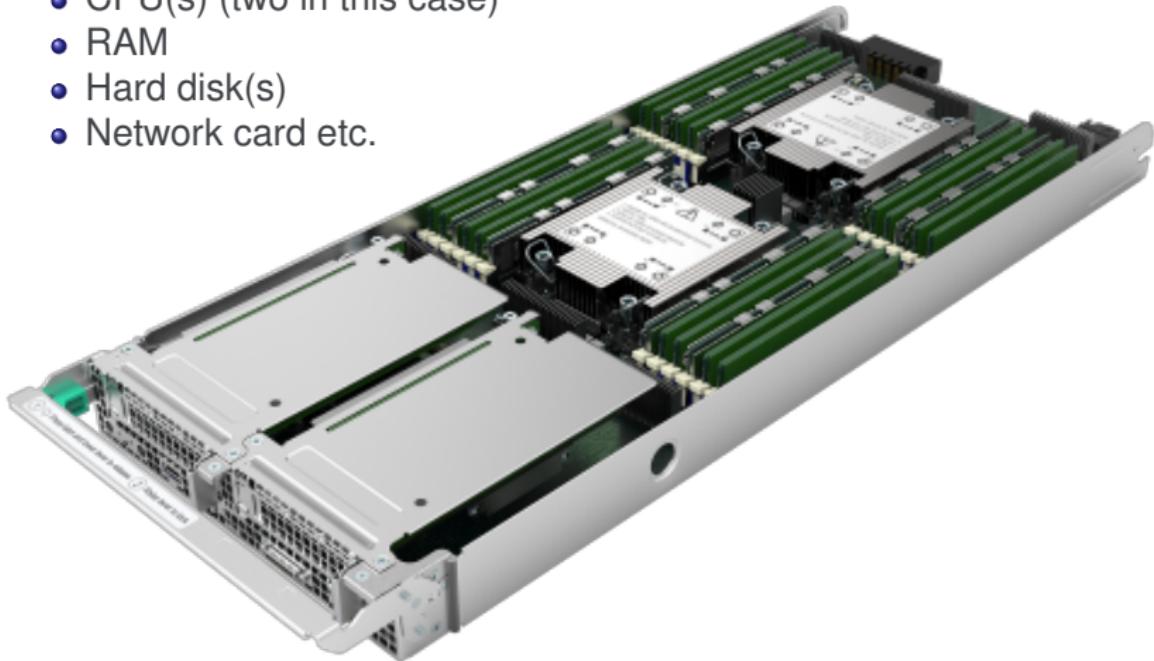
Server Rack

- Server racks hold servers, storage units, and switches

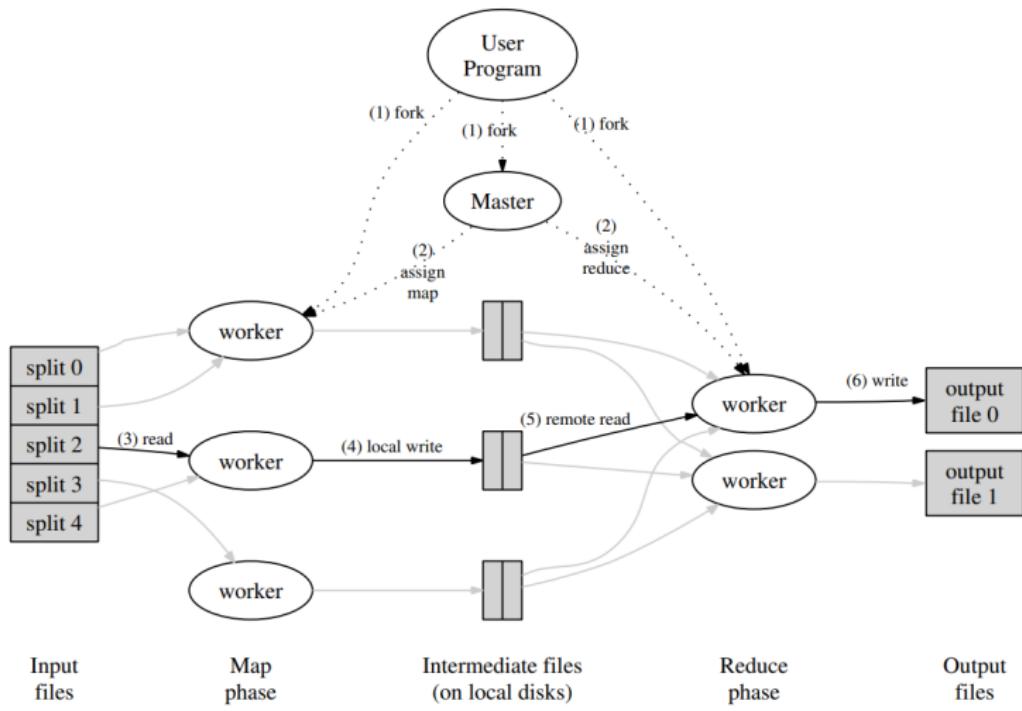


Rack-Mounted Server

- Contains inexpensive "commodity" hardware, similar to a desktop PC:
 - CPU(s) (two in this case)
 - RAM
 - Hard disk(s)
 - Network card etc.



The MapReduce Programming Model



Google File System

- Google's MapReduce is supported by a distributed file system called Google File System (GFS)
- When data is loaded onto the system it is divided into blocks
- Each block is stored on multiple different servers (typically 3)
- Files in GFS are write-once only, although data can be appended
- GFS is optimised for streaming reads of large files, rather than random access
- MapReduce tasks are divided into two phases:
 - Map tasks which are done concurrently on each block of data in a file, with each block being mapped over using the CPU of one of the servers on which it's stored
 - Reduce tasks then combine data to produce the final output

Advantages of Google File System

Advantages of Google's storage approach:

- It's resilient to server failure, because each block is stored on multiple servers
 - This is important, because failures will be frequent with thousands of servers in a cluster
- It allows parallel processing of the data, because the blocks of a single file are stored on different servers, each with its own CPU(s)
- It has huge capacity, because the file system uses hundreds or thousands of servers, each with hundreds of TBs of hard disk space

Advantages of Google File System

- It allows data sharing, because any user with access to GFS has access to the data, and can analyse it with MapReduce
- It's scalable – it scales linearly to handle larger data by adding more servers to the cluster
- Provides a clean abstraction for parallel programming:
 - MapReduce isn't essential for operating on big data, but it makes it easier
 - Google developed MapReduce to simplify big data tasks that engineers were previously implementing in more manual, *ad hoc*, ways.

Using MapReduce for Smaller Datasets

- In Google's original white paper, individual files were usually small enough to fit onto a single server (113GB vs 320GB)
 - However, the entire file store was far too large to fit onto a single server
- Big data middleware is useful for large files that can still technically fit onto a single server.
 - You still get the advantages of resilience, parallel execution of MapReduce, data sharing, etc.
- However, it's less well-suited for very small files:
 - GFS block size is 64MB, so files smaller than this will be stored as a single block, meaning that you'll get no parallel processing speed-up when using MapReduce

3

Distributed File Systems and MapReduce

- Google MapReduce and Google File System
- **Motivating Examples from Search Engine Design**
- Hadoop
- The Hadoop Distributed File System (HDFS)
- Orchestration of Hadoop MapReduce

Attendance: 500240

Motivating Examples

- In this part, we'll get an idea of just how useful MapReduce is for applying data processing to big data.

Word Counting

- Consider the problem of counting the number of occurrences of each word in a large collection of documents
 - This is an important step for the natural language processing that we'll do in week 9
 - We'll use these word counts to create document representations, and these will form the basis of our search engine designs

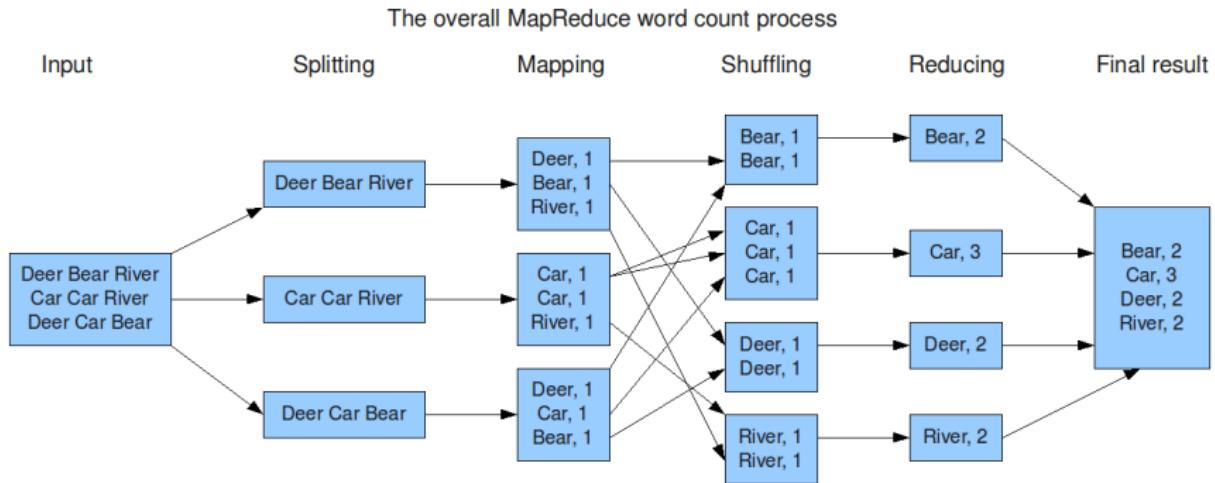
Word Counting

- MapReduce pseudo-code for this problem:

```
map(String key, String value):
    // key: document name
    // value: document contents
    for each word w in value:
        EmitIntermediate(w, "1");

reduce(String key, Iterator values):
    // key: a word
    // values: a list of counts
    int result = 0;
    for each v in values:
        result += ParseInt(v);
    Emit(AsString(result));
```

MapReduce: Word Count



Count of URL Access Frequency

- Dataset: a log of web page requests
- Output: a count of the number of requests for each web page
- MapReduce implementation:
 - The map function processes logs of web page requests and outputs (URL, 1) for each web request
 - The reduce function adds together all values for the same URL and emits a (URL, total count) pair

Reverse Web-Link Graph

- Dataset: a file of web pages and the URLs that each one links to
- Output: a mapping from each web page to the list of the web pages that link to it
- MapReduce implementation:
 - The map function outputs (target, source) pairs for each link to a target URL found in a page named source.
 - The reduce function concatenates the list of all source URLs associated with a given target URL and emits the pair: (target, list[source])

Term-Vector per Host

- Dataset: A collection of text files
- Output: A term vector summarises the most important words that occur in a document or a set of documents as a list of `<word, frequency>` pairs.
- MapReduce implementation:
 - The map function emits a `<hostname, term vector>` pair for each input document (where the `hostname` is extracted from the URL of the document).
 - The reduce function is passed all per-document term vectors for a given host. It adds these term vectors together, throwing away infrequent terms, and then emits a final `<hostname, term vector>` pair.

Inverted Index

- Dataset: A collection of text files
- Output: A mapping from each word to the list of documents that contain the word
- MapReduce implementation:
 - The map function parses each document, and emits a sequence of (word, document ID) pairs.
 - The reduce function accepts all pairs for a given word, sorts the corresponding document IDs and emits a (word, list(document ID)) pair.

Distributed Grep

- Dataset: a collection of text files
- Output: for a given regular expression, the set of lines that match it
- MapReduce implementation:
 - The map function emits a line if it matches a supplied pattern.
 - The reduce function is an identity function that just copies the supplied intermediate data to the output.

3

Distributed File Systems and MapReduce

- Google MapReduce and Google File System
- Motivating Examples from Search Engine Design
- **Hadoop**
- The Hadoop Distributed File System (HDFS)
- Orchestration of Hadoop MapReduce

Attendance: 500240

Hadoop

- Open source software framework, providing the functionality of Google's MapReduce and GFS
- Created by Doug Cutting and Mike Cafarella in 2005, based on Google's 2003 and 2004 white papers
- Cutting named the program after his son's toy elephant.



Who Uses Hadoop?



eHarmony®

facebook

IBM

The New York Times

JPMorganChase

twitter

intel

NETFLIX

rackspace
HOSTING

amazon.com

VISA

NING

SAMSUNG

YAHOO!

3

Distributed File Systems and MapReduce

- Google MapReduce and Google File System
- Motivating Examples from Search Engine Design
- Hadoop
- **The Hadoop Distributed File System (HDFS)**
- Orchestration of Hadoop MapReduce

Attendance: 500240

HDFS

- HDFS is Hadoop's implementation of the file system concept outlined by Google in their 2003 GFS white paper
- As with GFS:
 - HDFS works best with large files
 - Typically GBs per file
 - Files in HDFS are write-once only, but can be appended to
 - HDFS is optimised for streaming reads of large files and not random reads
- Files are split into 64MB blocks (for Hadoop v1) or 128MB blocks (for Hadoop v2)
- Each block is replicated across several servers (default 3)

HDFS Implementation Details

Hadoop clusters include the following daemons (i.e., programs) which work together to maintain **HDFS**:

- **NameNode:**

- Runs on the (unique) master server.
- Acts as the bookkeeper for HDFS.
- Keeps track of where blocks constituting copies of each file are stored.
- Monitors the health of the distributed file system.
- Is a single point of failure for the cluster (so it is commonly on a dedicated machine).

HDFS Implementation Details

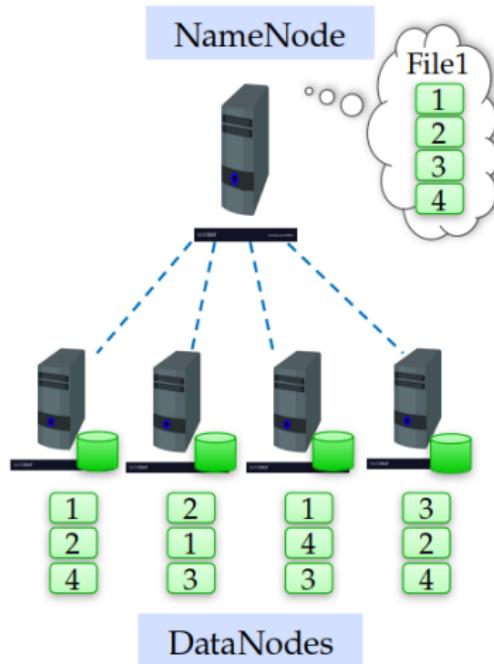
- **DataNodes:**

- Runs on each worker node.
- Stores the actual blocks (the client can upload files to DataNodes).
- Informs the NameNode of any local changes to files.
- Polls the NameNode for any required changes to blocks.

- **Secondary NameNode:**

- Takes snapshots of the NameNode's metadata.
- Purpose is to minimise the impact of NameNode failure.

HDFS Implementation Details



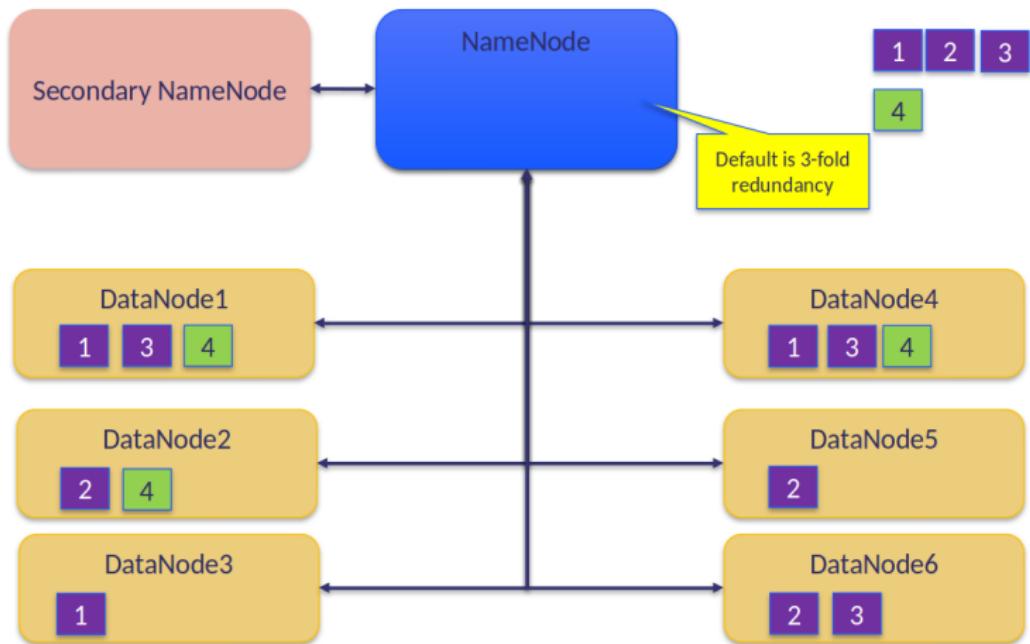
The NameNode

- Stores the [HDFS](#) file system information in the `FsImage`.
 - `FsImage` is a file stored on the OS filesystem that contains the complete directory structure (namespace) of the HDFS, with details about the location of the data on the Data Blocks and which blocks are stored on which node.
 - This file is used by the NameNode when it is started.
- Updates to the file system (add/remove blocks) do not change the `FsImage`.
 - They are instead written to a log file.
- When starting, the NameNode loads the `FsImage` file and then applies the changes in the log file.

The Secondary NameNode

- Not a backup for the NameNode
- Periodically reads the log file and applies the changes to the FsImage file, bringing it up to date
- Allows the NameNode to restart faster when required

HDFS Storage Example



Data Retrieval

- When a client wants to retrieve data:
 - Communicates with the **NameNode** to determine which blocks make up a file and on which data nodes those blocks are stored
 - Then communicates directly with the data nodes to read the data
 - So, no data passes through the master server

Hadoop MapReduce

- Written in Java
- Performs the MapReduce operation on HDFS data

The Mapper

- Reads data as key/value pairs
 - The key is often discarded
- Outputs zero or more key/value pairs
- E.g., for word count:

```
public static class TokenizerMapper
    extends Mapper<Object, Text, Text, IntWritable> {

    private final static IntWritable one = new IntWritable(1);
    private Text word = new Text();

    public void map(Object key, Text value, Context context
                    ) throws IOException, InterruptedException {
        StringTokenizer itr = new StringTokenizer(value.toString());
        while (itr.hasMoreTokens()) {
            word.set(itr.nextToken());
            context.write(word, one);
        }
    }
}
```

Shuffle and Sort

- Shuffle: copies the output from each Mapper using HTTP across the network.
- Sort: merge sorts Reducer inputs by key(s) (since different Mappers may have output the same key).
- As a result:
 - Output from the Mapper is sorted by key
 - All values with the same key go to the same worker server

The Reducer

- Called once for each unique key
- Gets a list of all values associated with a key as input
- The reducer outputs zero or more final key/value pairs
 - Usually just one output per input key
- E.g., for word count:

```
public static class IntSumReducer
    extends Reducer<Text, IntWritable, Text, IntWritable> {

    private IntWritable result = new IntWritable();

    public void reduce(Text key, Iterable<IntWritable> values,
                       Context context) throws IOException, InterruptedException {
        int sum = 0;
        for (IntWritable val : values) {
            sum += val.get();
        }
        result.set(sum);
        context.write(key, result);
    }
}
```

3

Distributed File Systems and MapReduce

- Google MapReduce and Google File System
- Motivating Examples from Search Engine Design
- Hadoop
- The Hadoop Distributed File System (HDFS)
- **Orchestration of Hadoop MapReduce**

Attendance: 500240

Orchestration of Hadoop MapReduce

- "Orchestration" refers to the process of implementing MapReduce computations on a Hadoop cluster
- Hadoop refers to each MapReduce call as a "job":
 - This "job" is broken down into multiple "tasks"
 - "Tasks" are then performed in parallel

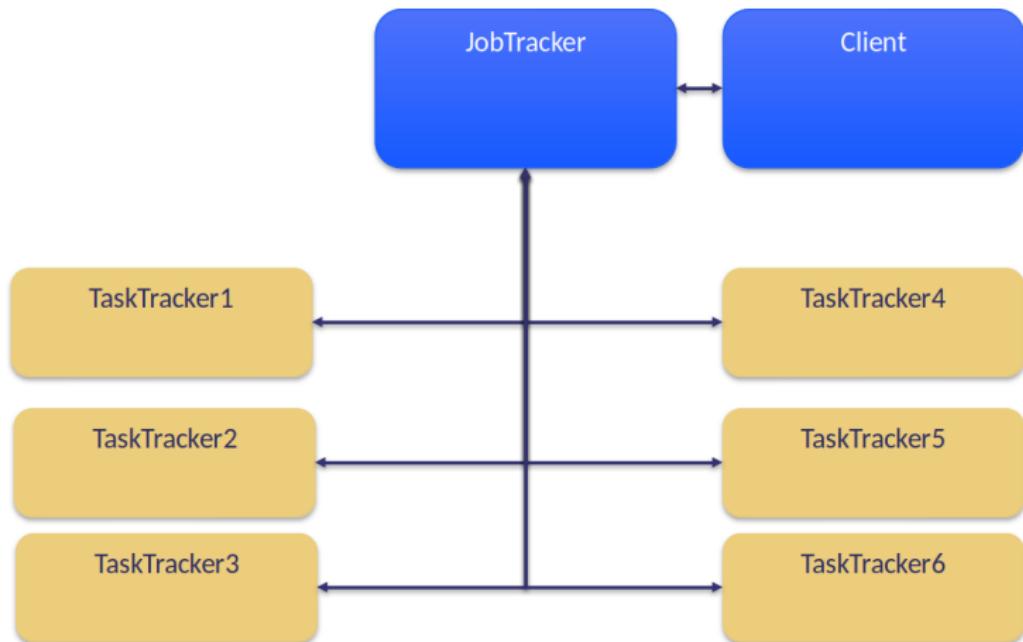
Orchestration of Hadoop MapReduce

- Hadoop clusters running **MapReduce v1** include the following daemons:
- **JobTracker:**
 - Runs on the single master server
 - The interface to clients wishing to run MapReduce jobs
 - Determines which computers do what and when
 - Will reassigned (map and reduce) “tasks” to computers if tasks fail
 - Including if it stops hearing the heartbeat signal
- **TaskTrackers:**
 - Run on each worker server
 - Spawns (multiple) Java Virtual Machines (JVMs) to execute tasks
 - Maintains a heartbeat with JobTracker

Orchestration of Hadoop MapReduce

- In MapReduce v2:
 - **JobTracker** was replaced by Resource Manager and Application Master
 - **TaskTracker** was replaced by Resource Manager, Application Manager, and Application Master

JobTracker and TaskTracker



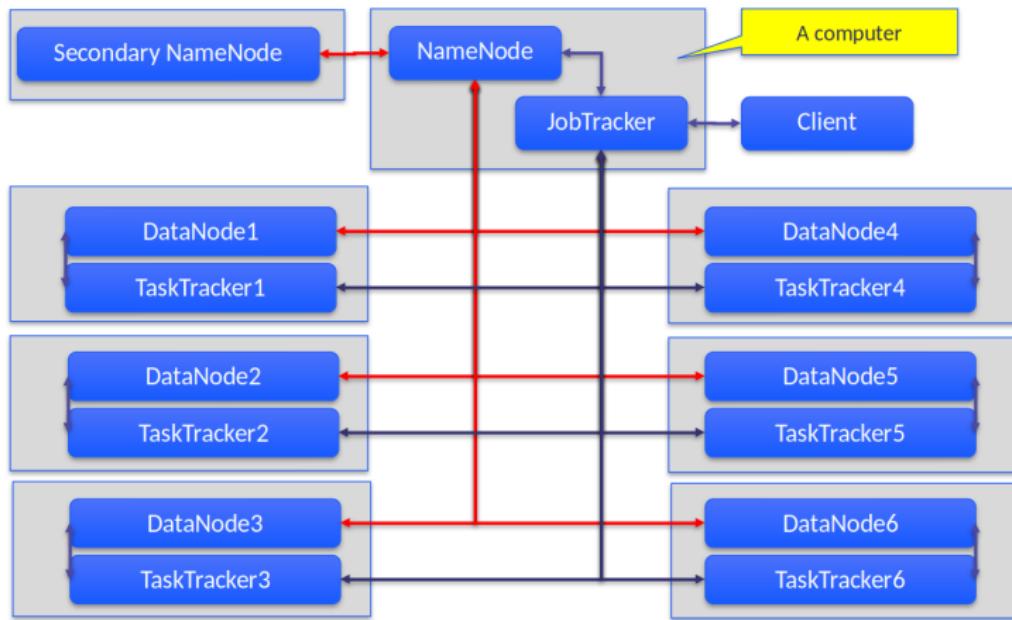
Execution Modes

- Hadoop can run in three execution modes:
 - 1 Fully distributed
 - 2 Pseudo-distributed
 - 3 Local

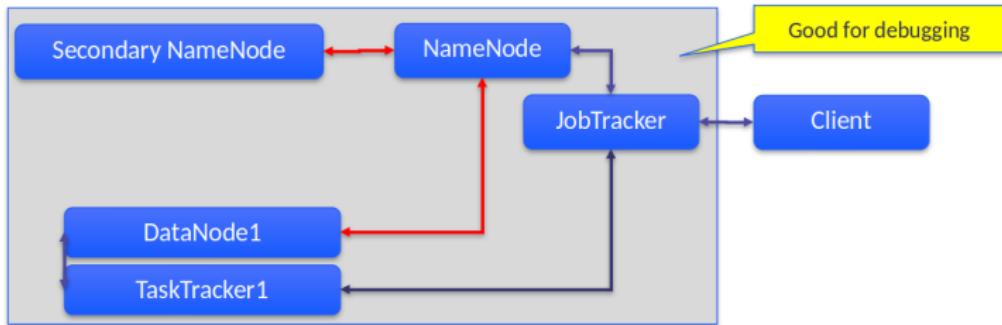
Fully Distributed Mode

- The following diagram represents a fully distributed mode in Hadoop.
 - The NameNode and Secondary NameNode handle metadata.
 - DataNodes store the actual data.
 - JobTracker and TaskTrackers manage the computational tasks.
- Key Components:
 - NameNode: Manages file system namespace and access.
 - DataNodes: Store data and report to NameNode.
 - JobTracker: Schedules and tracks jobs.
 - TaskTrackers: Execute tasks assigned by JobTracker.

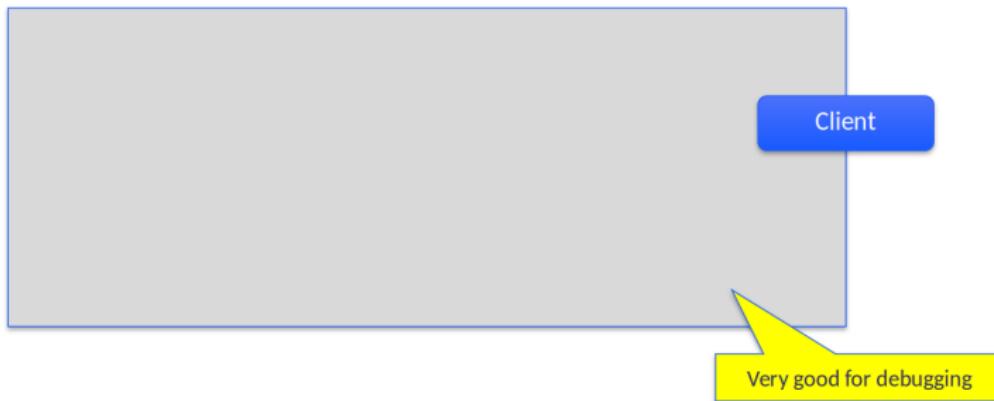
Fully Distributed Mode



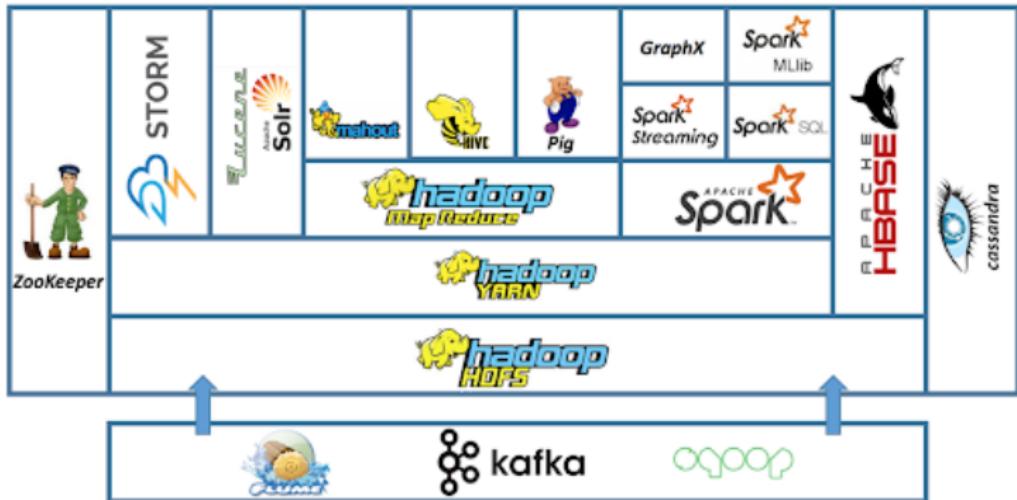
Pseudo-Distributed Mode



Local Mode



The Hadoop V2 Ecosystem



Spark

4

Spark

- Spark Programming Model
- Transformations
- Actions
- DataFrames
- Spark Computation
- Spark Architecture
- Spark Debugging and Running Modes

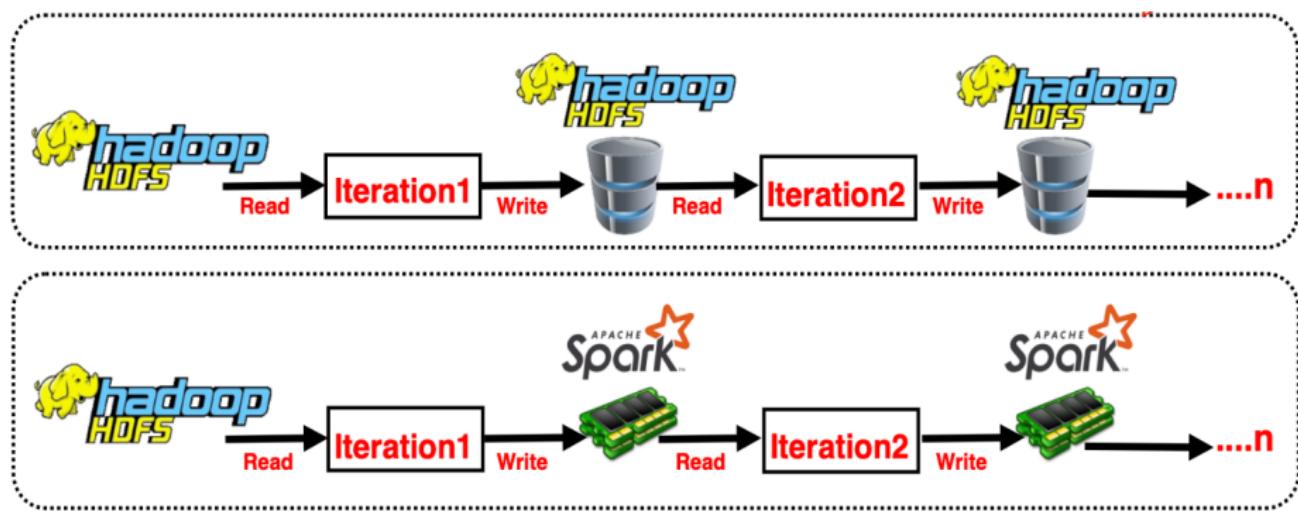
Attendance: 500240

Spark

- Developed as an improvement of Hadoop MapReduce.
- Works with any Hadoop-supported storage system (HDFS, S3, Avro, etc.), but not a modified version of Hadoop.
- Generally faster than Hadoop MapReduce:
 - Spark can perform in-memory processing, while Hadoop MapReduce has to read from/write to a disk (100x faster)
- Supports a broader set of computations than Hadoop
- APIs in Java, Scala, Python, and R (often 2-10x less code vs Hadoop)

Hadoop Vs Spark

- Spark keeps data in-memory while Hadoop MapReduce writes to disk after each map and reduce operation.
- This makes Spark 10x-100x faster than Hadoop for MapReduce computations



Word Count in Spark and Hadoop

Spark also has a concise API, often resulting in far less code than Hadoop:

```
1 public class WordCount {
2     public static class TokenizerMapper
3         extends Mapper<Object, Text, Text, IntWritable>{
4             private final static IntWritable one = new IntWritable();
5             private Text word = new Text();
6
7             public void map(Object key, Text value, Context context
8                             ) throws IOException, InterruptedException {
9                 StringTokenizer itr = new StringTokenizer(value.toString());
10                while (itr.hasMoreTokens()) {
11                    word.set(itr.nextToken());
12                    word.set(itr.nextToken());
13                    context.write(word, one);
14                }
15            }
16        }
17
18        public static class IntSumReducer
19            extends Reducer<Text,IntWritable,Text,IntWritable> {
20            private IntWritable result = new IntWritable();
21
22            public void reduce(Text key, Iterable<IntWritable> values,
23                               Context context
24                               ) throws IOException, InterruptedException {
25                int sum = 0;
26                for (IntWritable val : values) {
27                    sum += val.get();
28                }
29                result.set(sum);
30                context.write(key, result);
31            }
32        }
33
34        public static void main(String[] args) throws Exception {
35            Configuration conf = new Configuration();
36            String[] otherArgs = new GenericOptionsParser(conf, args).getRemainingArgs();
37            if (otherArgs.length < 2) {
38                System.out.println("Usage: wordcount <in> [<in>... <out>]");
39                System.exit(1);
40            }
41            Job job = new Job(conf, "word count");
42            job.setMapperClass(TokenizerMapper.class);
43            job.setMapperClass(TokenMapper.class);
44            job.setCombinerClass(IntSumReducer.class);
45            job.setReducerClass(IntSumReducer.class);
46            job.setOutputKeyClass(Text.class);
47            job.setOutputValueClass(IntWritable.class);
48            for (int i = 0; i < otherArgs.length - 1; i++) {
49                FileInputFormat.addInputPath(job, new Path(otherArgs[i]));
50            }
51            FileOutputFormat.setOutputPath(job,
52                new Path(otherArgs[otherArgs.length - 1]));
53            System.exit(job.waitForCompletion(true) ? 0 : 1);
54        }
55    }
```

```
1 val f = sc.textFile(inputPath)
2 val w = f.flatMap(l => l.split(" ")).map(word => (word, 1)).cache()
3 w.reduceByKey(_ + _).saveAsText(outputPath)
```

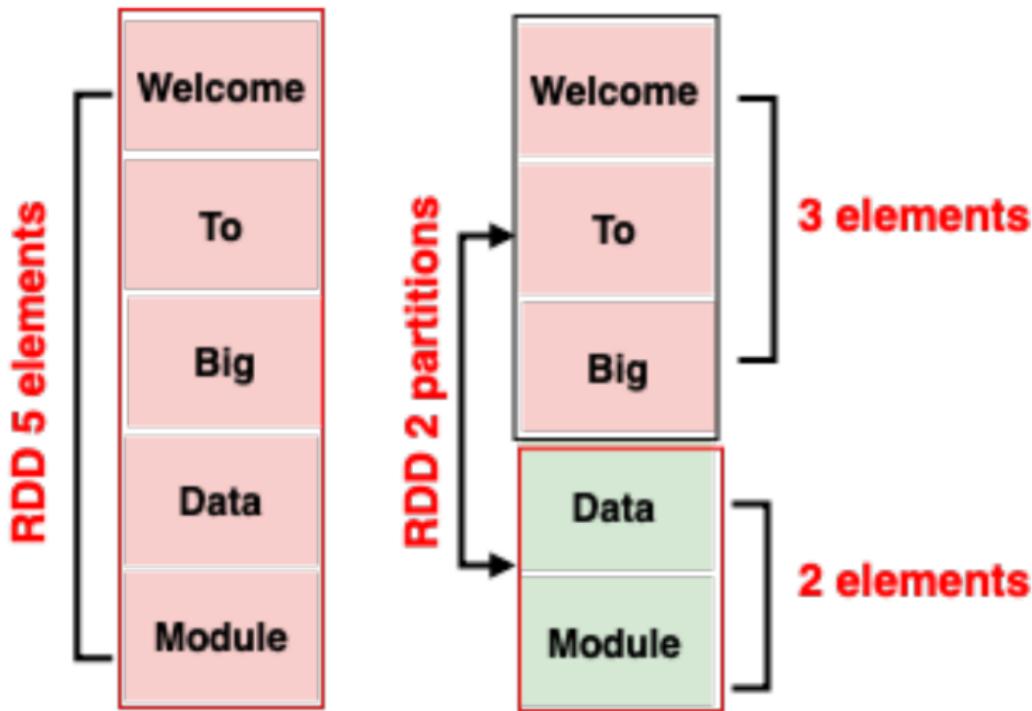
Figure: Word Count in Spark

Figure: Word Count in Hadoop

Resilient Distributed Datasets

- Data is represented in spark as **Resilient Distributed Datasets** (RDDs)
- RDDs are immutable collections of items/objects
- As we saw last week with HDFS, RDDs are split into partitions and can be stored on the worker nodes of the Spark cluster
- RDD operations can be performed on partitions in parallel
 - **Resilient** because each RDD partition is stored on multiple nodes of a Spark cluster
 - **Distributed** because different RDD partitions are spread across the cluster nodes, and can be processed in parallel

RDD Partitions



Immutability

- RDDs are **immutable** meaning that once they're created they can't be modified.
- This means that every RDD transformation creates a new RDD (just like lists in Haskell)

```
RDD1  
input_file = sc.textFile("/usr/local/spark/input.txt")  
  
RDD2  
map = input_file.flatMap(lambda line: line.split(" ")).map(lambda word: (word, 1))  
  
RDD3  
counts = map.reduceByKey(lambda a, b: a + b)  
  
RDD4  
counts.saveAsTextFile("/path/to/output/")
```

Creating RDDs

There are two ways to create an RDD:

- 1 Parallelising an existing collection in your driver program:



```
# Parallelize in Python
wordsRDD = sc.parallelize(["fish", "cats", "dogs"])
```



```
// Parallelize in Scala
val wordsRDD= sc.parallelize(List("fish", "cats", "dogs"))
```

```
// Parallelize in Java
JavaRDD<String> wordsRDD = sc.parallelize(Arrays.asList("fish", "cats", "dogs"));
```

- 2 Reading from an external storage system or filesystem:

```
# Read a local txt file in Python
input_file = sc.textFile("/usr/local/spark/input.txt")
```

Installing and Using PySpark

Spark is very easy to run locally with Python using PySpark:

- Install PySpark using pip:

```
$ pip install pyspark
```

- Import PySpark, initialize a SparkContext and run some commands:

```
$ python
>>> from pyspark import SparkContext
>>> sc = SparkContext("local", "PySpark Example")
>>> data = [1, 2, 3, 4, 5]
>>> rdd = sc.parallelize(data)
>>> squared_rdd = rdd.map(lambda x: x * x)
>>> result = squared_rdd.collect()
>>> print(result) # Output: [1, 4, 9, 16, 25]
[1, 4, 9, 16, 25]
>>> sc.stop()
```

The Spark Application UI

Visiting <http://localhost:4040/>, we see:

The screenshot shows the PySpark Example - Details for Job 0 - Chromium application UI. The browser title bar reads "PySpark Example - Details for Job 0 - Chromium". The address bar shows "localhost:4040/jobs/jobId". The main content area has a header "Details for Job 0" with sub-sections: "Status: SUCCEEDED", "Submitted: 2024/10/06 11:43:47", and "Duration: 0.6 s". Below this is a section titled "Completed Stages: 1" with a "Event Timeline" table. The table has columns for Executors (Added, Removed) and Stages (Completed, Failed, Active). A blue arrow points from the "Completed" stage row to the "Sun 6 October" date. The timeline shows a single event: "Executor driver added" at 11:34. Below the timeline is a "DAG Visualization" showing a single stage named "Stage 0" with one task labeled "parallelize". At the bottom, there are sections for "Completed Stages (1)" and "Stage Id 0" (with a "collect at <address>:1" entry), along with pagination controls.

PySpark Example - Details for Job 0 - Chromium

localhost:4040/jobs/jobId

PySpark Example application UI

Details for Job 0

Status: SUCCEEDED
Submitted: 2024/10/06 11:43:47
Duration: 0.6 s

Completed Stages: 1

Event Timeline

Executor driver added

11:34 11:35 11:36 11:37 11:38 11:39 11:40 11:41 11:42 11:43

Sun 6 October

DAG Visualization

```
graph TD; Stage0[Stage 0] --> Parallelize[parallelize]
```

Completed Stages (1)

Page: 1

| Stage Id | Description | Submitted | Duration | Tasks: Succeeded/Total | Input | Output | Shuffle Read | Shuffle Write |
|----------|------------------------|---------------------|----------|------------------------|-------|--------|--------------|---------------|
| 0 | collect at <address>:1 | 2024/10/06 11:43:47 | 0.4 s | 1/1 | | | | |

Page: 1

1 Pages. Jump to 1 Show 100 Items in a page Go

1 Pages. Jump to 1 Show 100 Items in a page Go

The Spark Stack

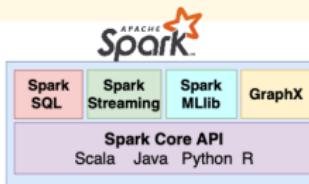
Spark Core, this is the heart of Spark, and is responsible for management functions such as task scheduling.

Spark SQL is a Spark module for structured data processing.

Spark Streaming is a framework for ingesting real-time streaming data from various sources.

Spark MLlib is a library of machine-learning algorithms (e.g., logistic regression, naïve Bayes classification, SVMs, decision trees, random forests, linear regression, and k-means clustering).

GraphX are data structures comprising vertices and the edges for graph and graph parallel computations.



Spark MLLib in Context

- A word of caution: Spark isn't always the best choice for a given machine learning application.
- Here are the alternatives:
 - **Spark MLLib:** Designed specifically to handle large-scale data processing across distributed computing environments. Ideal for big data applications.
 - **Python scikit-learn AKA sklearn:** A python library for machine learning. We'll use this in this course.
 - **R:** a language for statistics and machine learning.
 - **TensorFlow/PyTorch:** frameworks for designing and training deep neural networks.

Comparison of Machine Learning Frameworks

| Framework | Dataset Size | ML Techniques |
|------------|----------------------|-------------------------------|
| Spark | Large-scale | Basic |
| sklearn | Small to Medium | Basic, Intermediate |
| R | Small to Medium | Basic, Intermediate, Advanced |
| TF/PyTorch | Small to Large-scale | Deep Learning |

- R is the gold standard for stats and machine learning on small to medium sized datasets (with the exception that deep learning is mostly done in Python with TensorFlow or PyTorch)
- We don't use it in this course because there's a (small) learning curve due to the fact that it provides its own programming language
- However, the language is very easy to learn, so R is ideal for anything most stats and ML projects

4

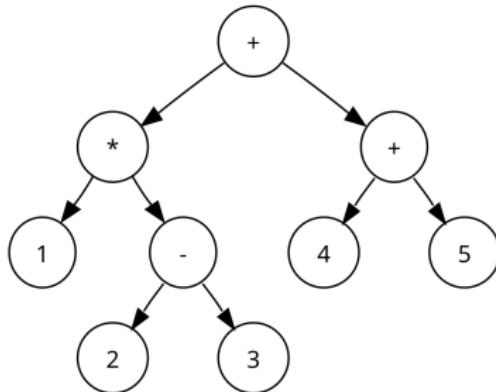
Spark

- Spark Programming Model
- Transformations
- Actions
- DataFrames
- Spark Computation
- Spark Architecture
- Spark Debugging and Running Modes

Attendance: 500240

Directed Acyclic Graphs and Data Flows

- Spark programming model is based on parallelizable operators, which can be composed together to form a **directed, acyclic graph (DAG)**
 - This is a directed graph with no directed cycles
 - E.g. here's a DAG for the trivial computation $1 * (2 - 3) + 4 + 5$:



- A **data flow** is composed of any number of data sources, operators, and data sinks. This is done by connecting their inputs and outputs.

DAGs in Spark and Lazy Evaluation

- In Spark, DAGs are used to represent a sequence of computations performed on data, where
 - Each **node** in the graph is an RDD
 - Each **edge** in the graph is a transformation that creates a new RDD from an existing one
- The DAGs are evaluated **lazily** (just like functions in Haskell)
- This allows Spark to optimise the execution plan by reordering and combining transformations.
- This optimisation leads to efficient in-memory execution, which **explains why Spark is so much faster than Hadoop**

Directed Acyclic Graphs in Spark



Spark RDD Operations

Spark RDD support two types of operations: transformations and actions.

- i. **Transformations** can be applied to an RDD to create a new RDD
 - These are the DAG edges
 - E.g. `x.map(f)` creates a new RDD by applying `f` to each element in `x`.
- ii. **Actions** compute a result based on an RDD and either return a value or save it to an external storage system (e.g., HDFS).
 - E.g. `x.saveAsTextFile(path)` writes one RDD element per line to `path`

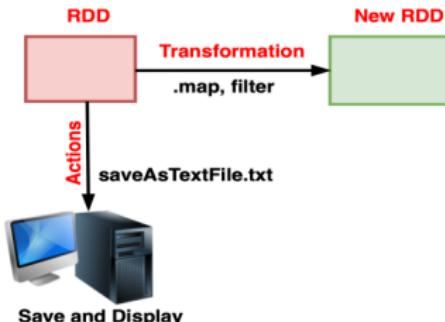
Spark RDD Operations

Transformation

```
input_file = sc.textFile("/usr/local/spark/input.txt") map  
= input_file.flatMap(lambda line: line.split(" "  
)).map(lambda word: (word, 1))  
  
counts = map.reduceByKey(lambda a, b: a + b)
```

```
counts.saveAsTextFile("/path/to/output/")
```

Action

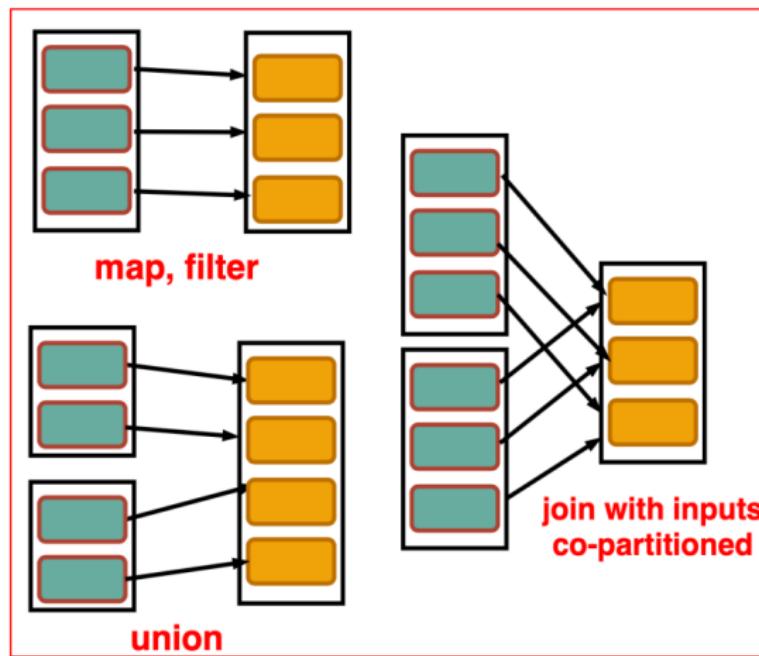


Narrow Vs Wide Transformation

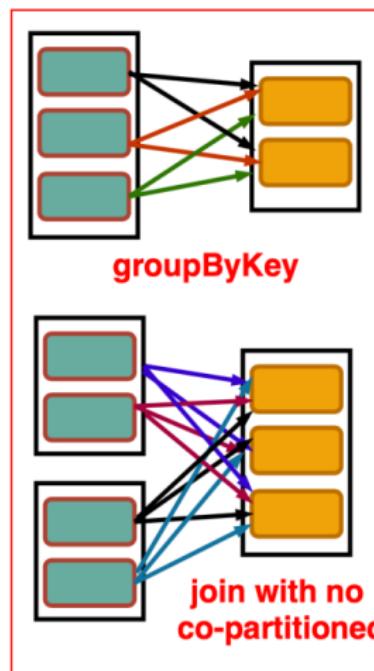
Transformations can be narrow or wide:

- **Narrow Transformation**: each partition of the parent RDD is used by at most one partition of the child RDD
- **Wide (AKA shuffle) transformations**: multiple child RDD partitions may depend on a single parent RDD partition

Narrow Transformations



Wide Transformations



RDD Operations in Spark

| | General | Math / Statistical | Set Theory / Relational | Data Structure / I/O |
|-----------------|---|---|--|--|
| Transformations | <ul style="list-style-type: none">mapfilterflatMapmapPartitionsmapPartitionsWithIndexgroupByKeysortBy | <ul style="list-style-type: none">samplerandomSplit | <ul style="list-style-type: none">unionintersectionsubtractdistinctcartesianzip | <ul style="list-style-type: none">keyByzipWithIndexzipWithUniqueIdzipPartitionscoalescerepartitionrepartitionAndSortWithinPartitionspipe |
| Actions | <ul style="list-style-type: none">reducecollectaggregatefoldfirsttakeforeachtoptreeAggregatetreeReduceforeachPartitioncollectAsMap | <ul style="list-style-type: none">counttakeSamplemaxminsumhistogrammeanvariancestdevsampleVariancecountApproxcountApproxDistinct | <ul style="list-style-type: none">takeOrdered | <ul style="list-style-type: none">saveAsTextFilesaveAsSequenceFilesaveAsObjectFilesaveAsHadoopDatasetsaveAsHadoopFilesaveAsNewAPIHadoopDatasetsaveAsNewAPIHadoopFile |

Highlighted operations are hard

4

Spark

- Spark Programming Model
- **Transformations**
- Actions
- DataFrames
- Spark Computation
- Spark Architecture
- Spark Debugging and Running Modes

Attendance: 500240

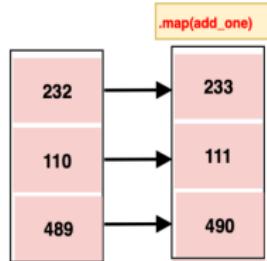
map and flatMap

- **map:** takes as input a function which is applied to each element of the RDD and map each input item to another item

```
# transforming each element through a function
```

```
nums = sc.parallelize([232, 110, 489])  
big_no = nums.map(lambda x: [n + 1])
```

If we want to map
every record 1,
.map(lambda x: [x,1])



(232, 1), (110,1),
(489,1)

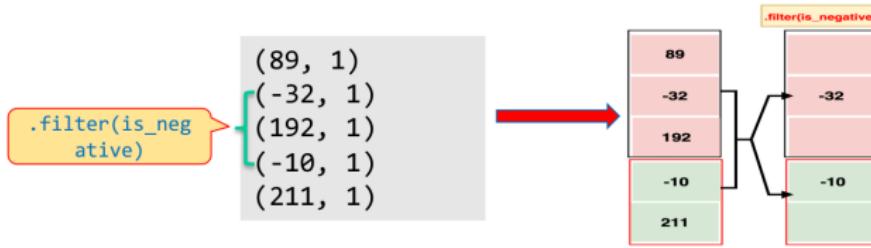
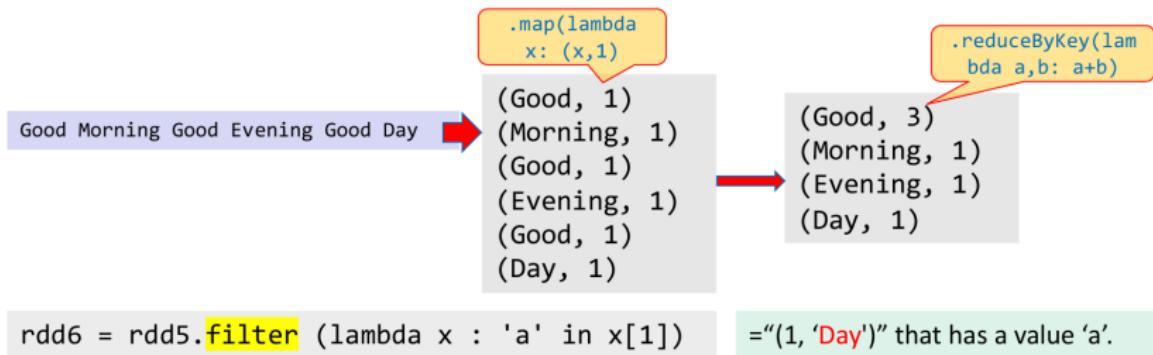
- **flatMap:** generates a new RDD by filtering the source dataset using the specified function

Good Morning Good Evening Good Day
Happy Birthday Happy New Year

```
lines.flatMap(lambda line: line.split()).collect()  
['Good', 'Morning', 'Good', 'Evening', 'Good', 'Day', 'Happy',  
'Birthday', 'Happy', 'New', 'Year']
```

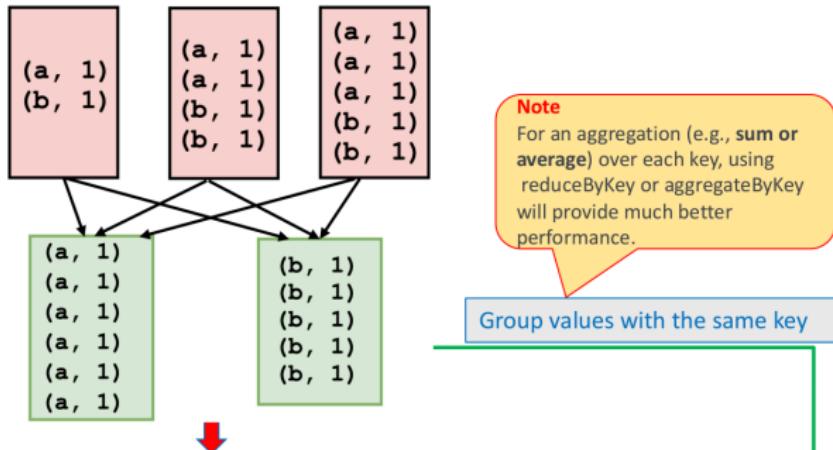
filter

filter generates a new RDD by filtering the source dataset using the specified function:



groupByKey

groupByKey groups the data based on a key and returns an RDD of grouped data:

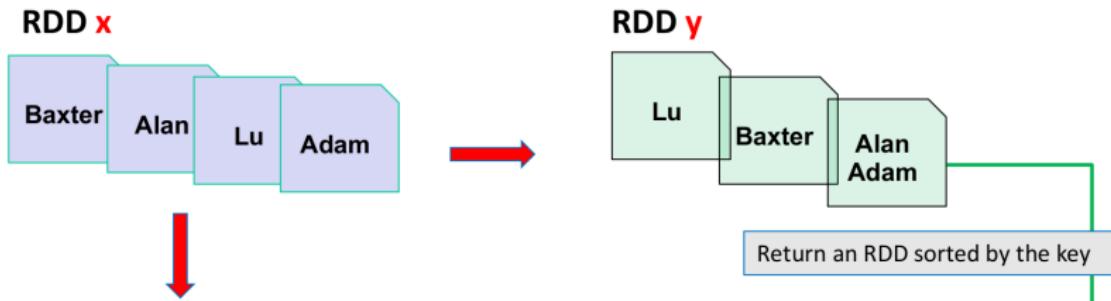


```
rdd = sc.parallelize([('a', 1), ('b', 1), ('a', 1), ('a', 1), ('b', 1), ('b', 1), ('a', 1), ('a', 1), ('a', 1), ('b', 1), ('b', 1)])
```

```
sorted(rdd.groupByKey().mapValues(list).collect()) [('a', [1, 1, 1, 1, 1, 1]), ('b', [1, 1, 1, 1, 1])]
```

groupBy

groupBy groups the data based on a user-provided function that is applied to each element in the RDD:



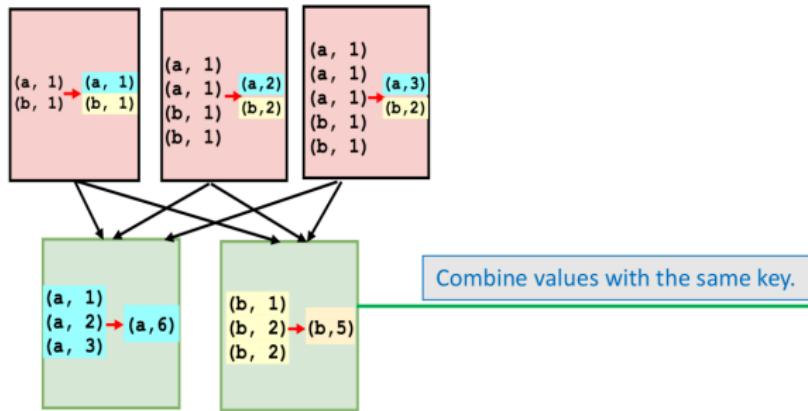
```
x = sc.parallelize(['Baxter', 'Alan', 'Lu', 'Adam'])
y = x.groupBy(lambda w: w[0])
```

```
print [(k, list(v)) for (k, v) in y.collect()]
```

```
x: ['Baxter', 'Alan', 'Lu', 'Adam']
y: [('A', ['Alan', 'Adam']), ('B', ['Baxter']), ('L', ['Lu'])]
```

reduceByKey

reduceByKey aggregates elements of the RDD using a specified associative and commutative function:

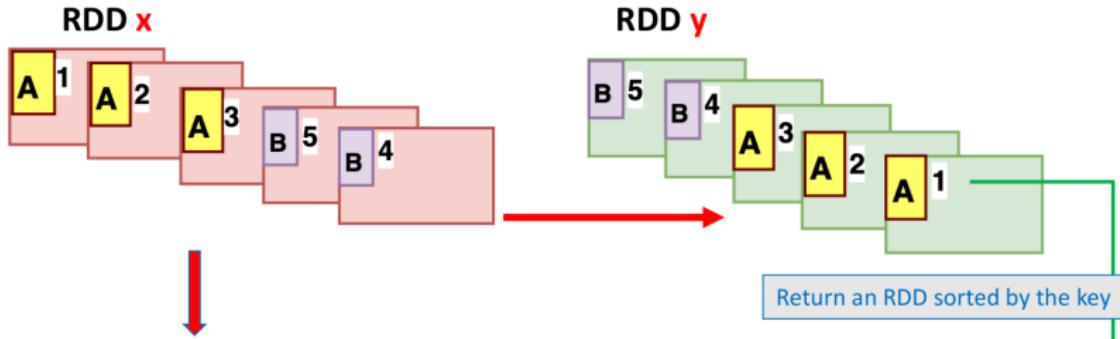


```
from operator import add
rdd = sc.parallelize([('a', 1), ('b', 1), ('a', 1), ('a', 1), ('b', 1),
                     ('b', 1), ('a', 1), ('a', 1), ('a', 1), ('b', 1), ('b', 1)])
```

```
sorted(rdd.reduceByKey(add).collect()) [('a', 6), ('b', 5)]
```

sortByKey

sortByKey sorts elements by key:

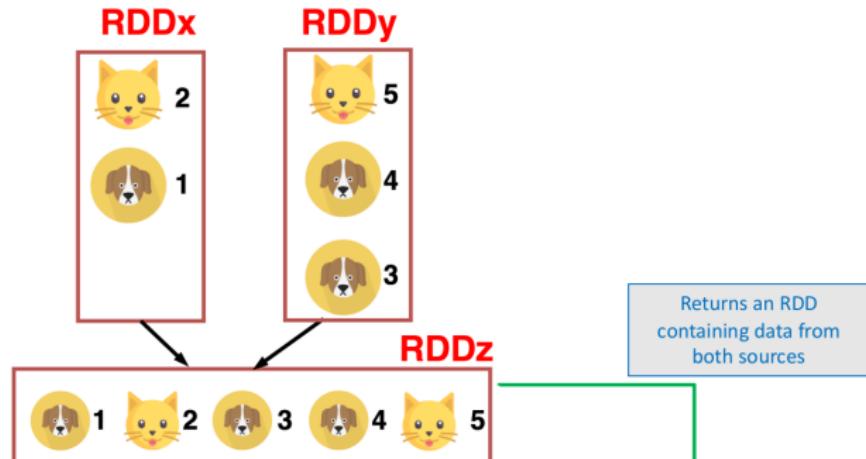


```
x = sc.parallelize([('B',4),('B',5),('A',3),('A',2),('A',1)])
y = x.sortByKey()
print(x.collect())
print(list((j[0], list(j[1])) for j in y.collect()))
```

```
x: [('B', 4), ('B', 5), ('A', 3), ('A', 2), ('A', 1)]
y: [('A', 1), ('A', 2), ('A', 3), ('B', 4), ('B', 5)]
```

union

union returns the union of two RDDs:

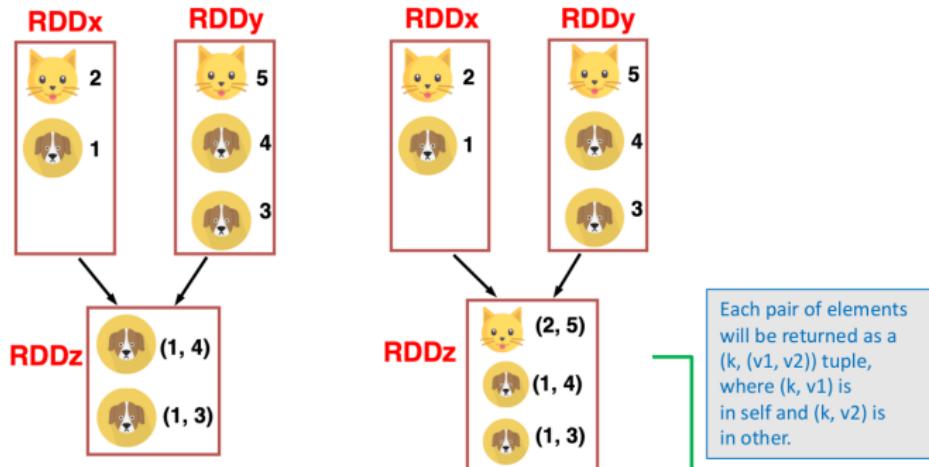


```
x = sc.parallelize([("dog",1), ("cat",2)])
y = sc.parallelize([("dog",3), ("dog",4), ("cat", 5)])
```

```
z = [("dog",1), ("cat",2), ("dog",3), ("dog",4), ("cat", 5)]
```

join

join returns the inner join of two RDDs:



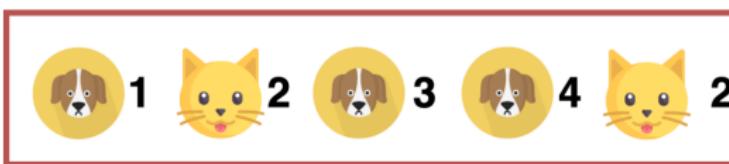
```
x = sc.parallelize([("dog", 1), ("cat", 2)])
y = sc.parallelize([("dog", 3), ("dog", 4), ("cat", 5)])
sorted(x.join(y).collect())
```

```
z: [('dog', (1,3)), ('dog', (1,4)), ('cat', (2, 5))]
```

distinct

distinct returns an RDD with distinct elements:

RDD x



Return a new RDD containing distinct value/s from the original RDD (omitting all duplicates)

```
x = sc.parallelize([("dog", 1), ("cat", 2), ("dog", 3), ("dog", 4),  
                   ("cat", 2)])  
y = x.distinct()  
print(y.collect())
```

```
z = [("dog", 1), ("cat", 2), ("dog", 3), ("dog", 4)]
```

4

Spark

- Spark Programming Model
- Transformations
- **Actions**
- DataFrames
- Spark Computation
- Spark Architecture
- Spark Debugging and Running Modes

Attendance: 500240

Actions

- **reduce** AKA **aggregate** applies a binary operator to an initial state and all elements in an RDD, and reduces this to a single state. The final state is converted into the final result by applying a finish function.
 - Note that `reduceByKey` is a transformation because it returns an RDD of keys and transformed values.
 - However, `reduce` is an action because it returns a single value rather than an RDD
- **fold** is similar to `reduce`. It aggregates the elements of each partition, and then the results for all the partitions. In contrast to `reduce`, it uses a neutral “zero value”.
- **collect**: Returns all the elements of the dataset as a list (NB: avoid using with large datasets because it will dump the entire dataset into your local computer’s memory).

Actions

- **first** returns the first element in this RDD.
- **take(n)** returns the first n elements from the dataset.
- **foreach(func)** applies a function to all elements of the RDD, generally used for side effects (e.g., saving data to a database). Similar to map, except that it is an action, so it does not produce a new RDD.
- **top(n)** returns the top n elements of the RDD, based on the default or specified ordering.

Actions

- **count** return the number of elements in this RDD.
- **takeSample(withReplacement, num, [seed])** returns a fixed-size sampled subset of this RDD.
- **max** returns the maximum element in the RDD (must be comparable).
- **min** returns the minimum element in the RDD (must be comparable).
- **sum** computes the sum of elements in the RDD.
- **histogram(buckets)** computes a histogram of the elements in the RDD, dividing them into buckets.
- **mean** returns the mean of the dataset.
- **variance** returns the variance of the dataset.
- **stdev** returns the standard deviation of the dataset.

4

Spark

- Spark Programming Model
- Transformations
- Actions
- **DataFrames**
- Spark Computation
- Spark Architecture
- Spark Debugging and Running Modes

Attendance: 500240

DataFrames

DataFrames are distributed collections of rows under named columns.

- Introduced to Spark in 2015
- Similar to Excel tables or Python DataFrames
- Build on top of RDDs
- Immutable in Nature: You can create a DataFrame, but you will not be able to change it. A DataFrame can be transformed just like an RDD.
- Lazy Evaluations: A task is not executed until an action is performed.
- Distributed: DataFrames, just like RDDs, are distributed in nature.

Creating a DataFrame in PySpark

```
>>> from pyspark.sql import SparkSession
>>> spark = SparkSession.builder \
...     .appName("PySpark Example") \
...     .getOrCreate()
>>> data = [
...     ("Alice", 25, "Data Scientist"),
...     ("Bob", 30, "Software Engineer"),
...     ("Catherine", 29, "Doctor")
... ]
>>> columns = ["Name", "Age", "Occupation"]
>>> df = spark.createDataFrame(data, columns)
>>> df.show()
+-----+---+-----+
|    Name| Age| Occupation|
+-----+---+-----+
|    Alice| 25| Data Scientist|
|      Bob| 30| Software Engineer|
| Catherine| 29|          Doctor|
+-----+---+-----+
```

How to Create DataFrames?

- **Using different data formats:** Data can be loaded from JSON, CSV, RDBMS, XML, or Parquet.
- **Loading from an existing RDD:** A DataFrame can be created from an already existing RDD.
- **Programmatically specifying a schema:** You can specify the schema programmatically to create a DataFrame.

4

Spark

- Spark Programming Model
- Transformations
- Actions
- DataFrames
- **Spark Computation**
- Spark Architecture
- Spark Debugging and Running Modes

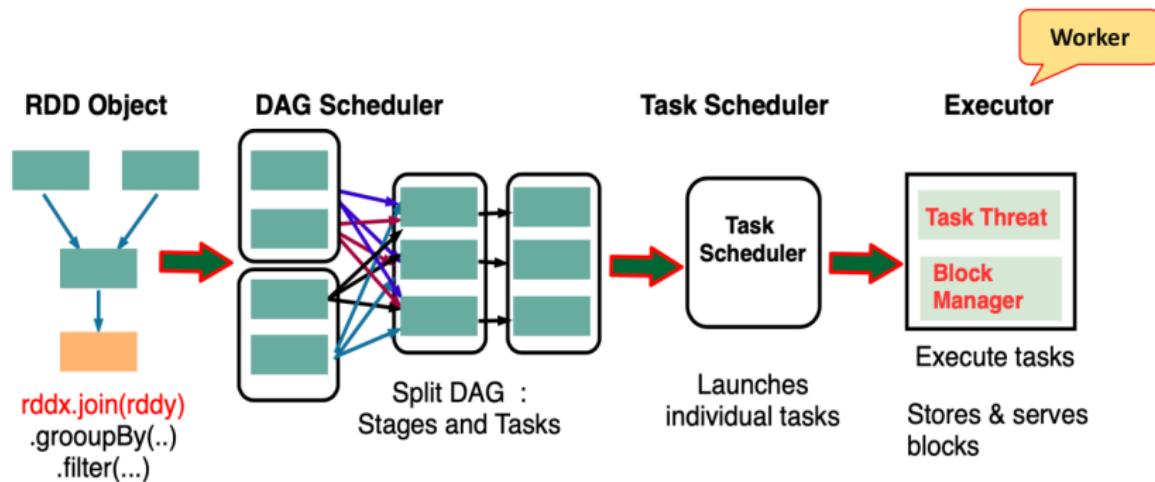
Attendance: 500240

Spark Components

- **Task:** Individual unit of work sent to one executor over a sequence of partitions.
- **Job:** Set of tasks executed as a result of an action.
- **Stage:** Set of tasks in a job that can be executed in parallel at the partition level.
- **RDD:** Parallel dataset with partitions.
- **DAG:** Logical graph of RDD operations.

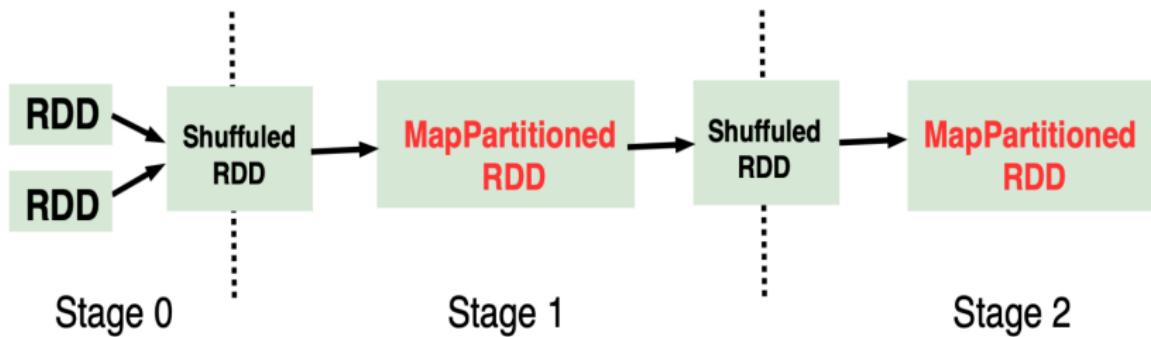
Job Scheduler

The job scheduler transforms a logical execution plan (i.e., RDD lineage of dependencies built using RDD transformations) to a physical execution plan (using stages).



Abstraction 2: RDD Stages

- Each **Spark stage** will dynamically allocate executors.
- Stages contain many neighbouring (in the lineage graph) transformations and actions.

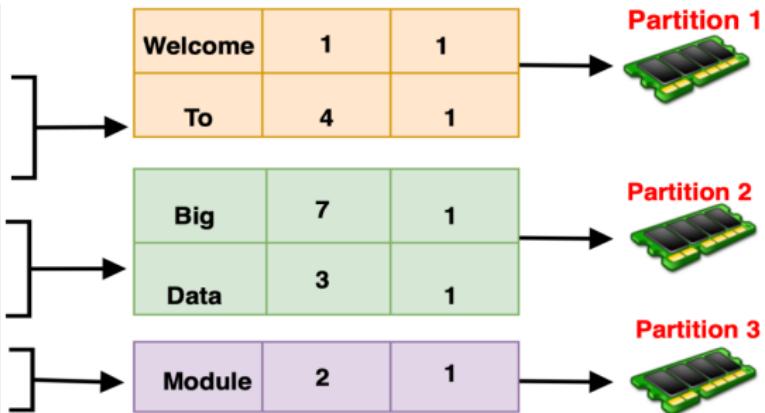


`rdd1.join(rdd2).groupByKey().filter()`

RDD

RDD

| Words | Index | Count |
|---------|-------|-------|
| Welcome | 1 | 1 |
| To | 4 | 1 |
| Big | 7 | 1 |
| Data | 3 | 1 |
| Module | 2 | 1 |



4

Spark

- Spark Programming Model
- Transformations
- Actions
- DataFrames
- Spark Computation
- **Spark Architecture**
- Spark Debugging and Running Modes

Attendance: 500240

Spark Architecture

Spark architecture is made up of a **driver** program and a set of **executor** programs.

Driver:

- Runs the `main()` function, located on a node in the cluster.
- Responsible for:
 - Maintaining information about the Spark applications.
 - Responding to a user's program or input, analyzing, distributing.
 - Scheduling work across the executors.

Executors:

- Perform all the data processing as assigned by the driver.
- Read from and write data to external sources.
- Store the computation results in-memory, cache, or on hard disk drives.
- Interact with the storage systems.

4

Spark

- Spark Programming Model
- Transformations
- Actions
- DataFrames
- Spark Computation
- Spark Architecture
- **Spark Debugging and Running Modes**

Attendance: 500240

Spark: Standalone Cluster

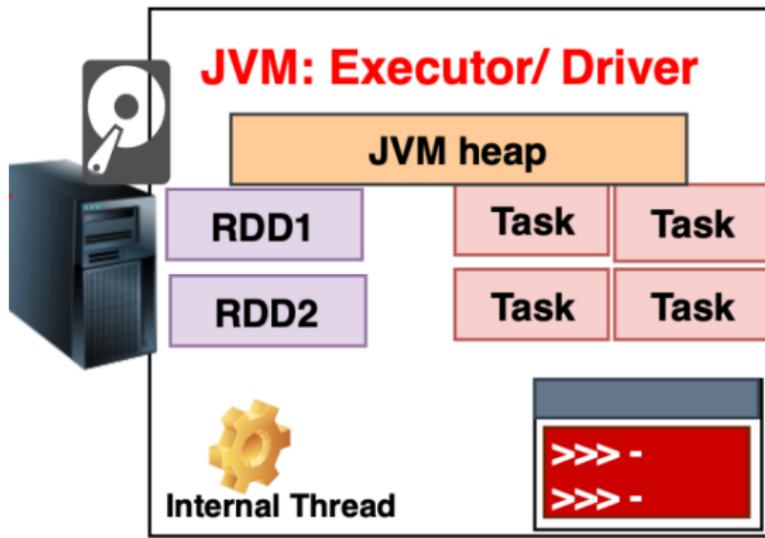
- Spark **local mode** and Spark **local cluster mode** are special cases of running Spark on a single machine.
 - **Local mode:** Runs Spark with one worker thread.
 - **Local[N]:** Runs Spark with N worker threads (like a thread pool).
 - **Local[*]:** Runs Spark with as many worker threads as there are logical cores on your machine.
- This mode is useful for running Spark applications on a local machine for testing or small workloads.

SparkContext Example in Local Mode

```
import pyspark
sc = pyspark.SparkContext('local[4]')
txt = sc.textFile('file:///usr/share/doc/python/copyright')
python_lines = txt.filter(lambda line: 'python' in line.lower())
print(python_lines.count())
```

- The SparkContext is created in local mode with 4 worker threads.
- The `textFile` method reads a file.
- The `filter` operation filters the lines containing the word "python."
- The `count` operation returns the number of matching lines.

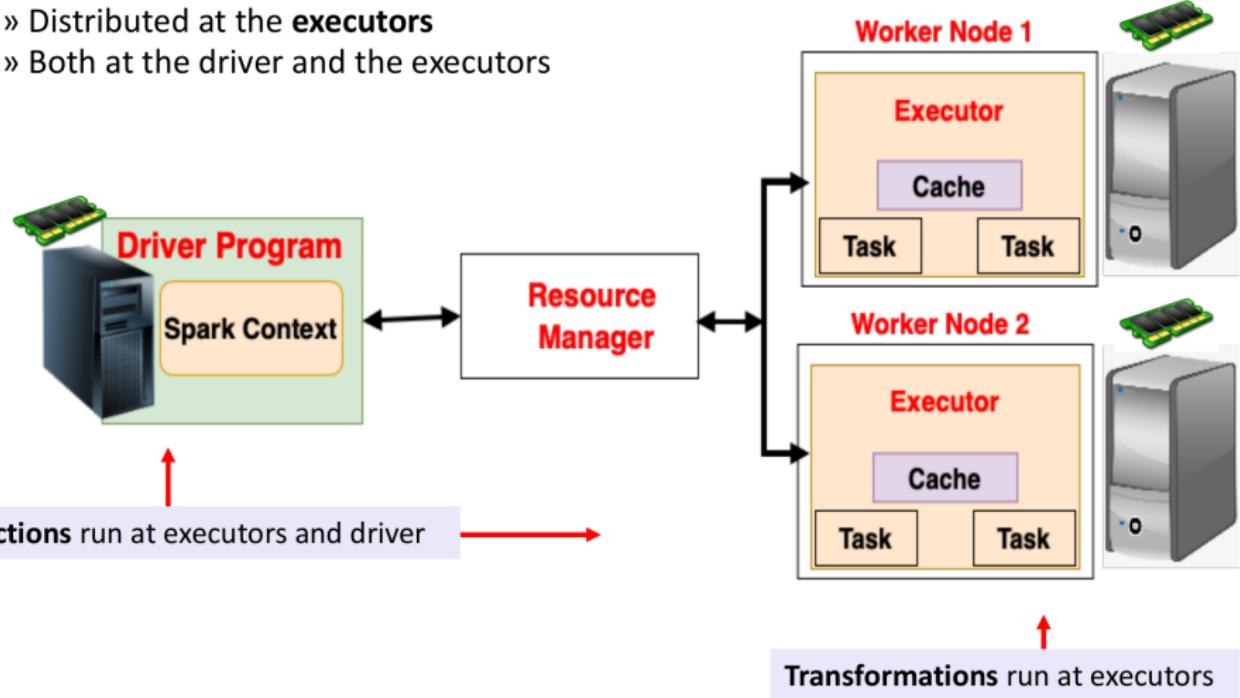
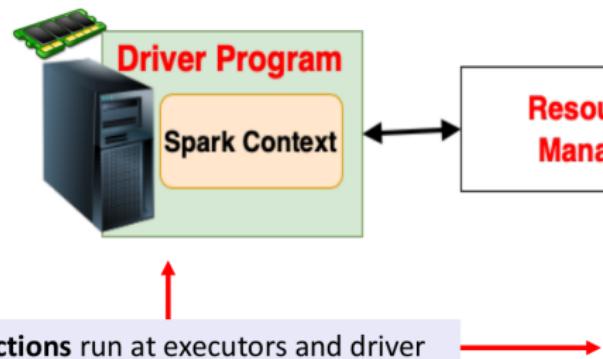
Spark Standalone Cluster



Cluster Mode

Where does code run?

- » Distributed at the **executors**
- » Both at the driver and the executors



References

- Feng, W. (2019). Learning Apache Spark with Python
- <https://databricks.com/glossary/spark-api>
- <https://realpython.com/pyspark-intro/>
- <https://www.cs.princeton.edu/courses/archive/fall13/cos518/papers/spark.pdf>

Classification

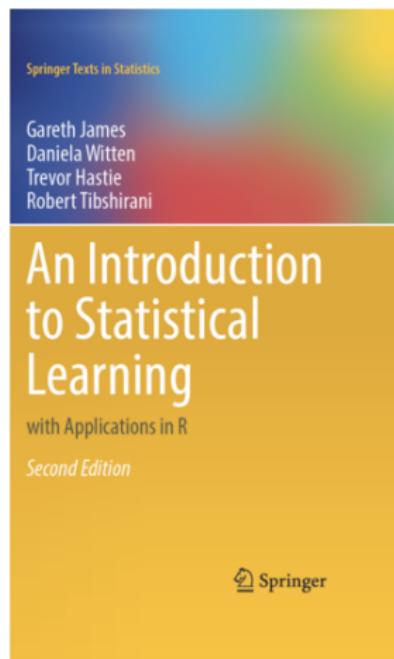
5 Classification

- Maths for Classifiers

Attendance: 500240

Recommended Reading for Classification

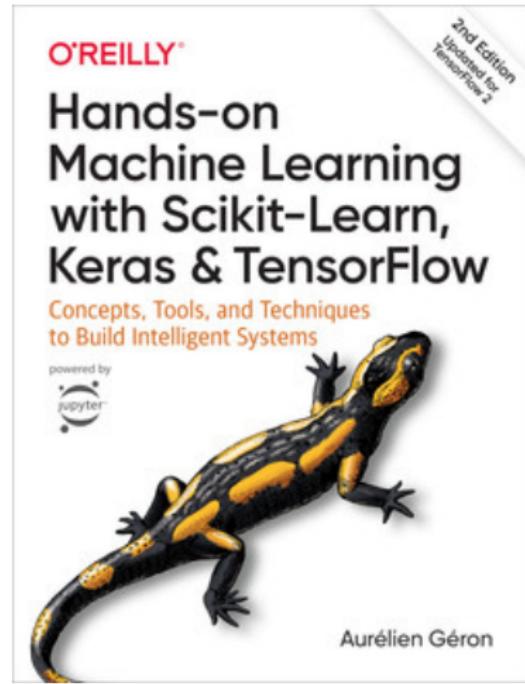
- This lecture is self-contained, but the following might also be useful:
- **An Introduction to Statistical Learning**
 - See chapter 9 for the maths in today's lecture
 - Examples in R (not used in this course)
 - Free online:
<https://www.statlearning.com/>
 - Accompanying lectures: <https://web.stanford.edu/~hastie/lectures.htm>



Recommended Reading for Classification

- **Hands-on Machine Learning, Geron, 3rd Edition 2022**

- Available online from the university library: <https://libguides.liverpool.ac.uk/online>
- Chapter 5 introduces classifiers without maths
- Examples in Python
- Code for today's lecture came from Hands-on ML
- Note: there's another book by the same name with a different author - don't pick that one up!



Supervised Versus Unsupervised Learning

- You've seen a few machine learning algorithms so far: PCA, clustering and linear regression.
- How do these fit into the broader space of ML algorithms?
- These are three types of machine learning algorithms: unsupervised learning, supervised learning and reinforcement learning
- **Unsupervised learning:**
 - The training data is a design matrix X , representing observations of a system
 - The aim is to find structure in this data
 - Typically, we're *compressing this design matrix* (making it narrower) in one way or another.
 - Examples:
 - PCA (week 5)
 - Clustering (week 6)

Supervised Versus Unsupervised Learning II

- **Supervised learning:**

- The training data is a design matrix X and a vector (or matrix) of responses y
- The aim is to train a model to predict the responses in y given the observations in X
- This model is then used to predict unseen responses from new observations
- Examples:
 - Linear regression (week 7)
 - K-nearest neighbours (today)
 - Linear support vector classifiers (today)
 - Support vector machines (today)

- **Reinforcement learning** is outside the scope of this course, but it's about learning to play games

Supervised Learning

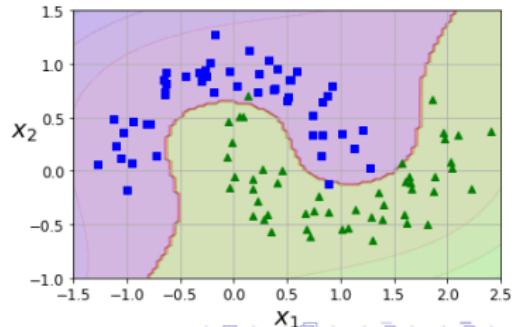
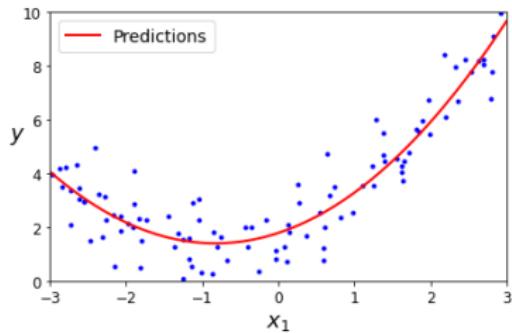
There are two types of supervised learning:

- **Regression:**

- The responses are real numbers or vectors thereof.
- Regression model outputs a real value or vector.
- E.g., linear regression (week 7).

- **Classification:**

- The responses are integers representing membership in a class.
- Classification model outputs a class.
- E.g., K-Nearest Neighbours, SVMs, linear SVCs.

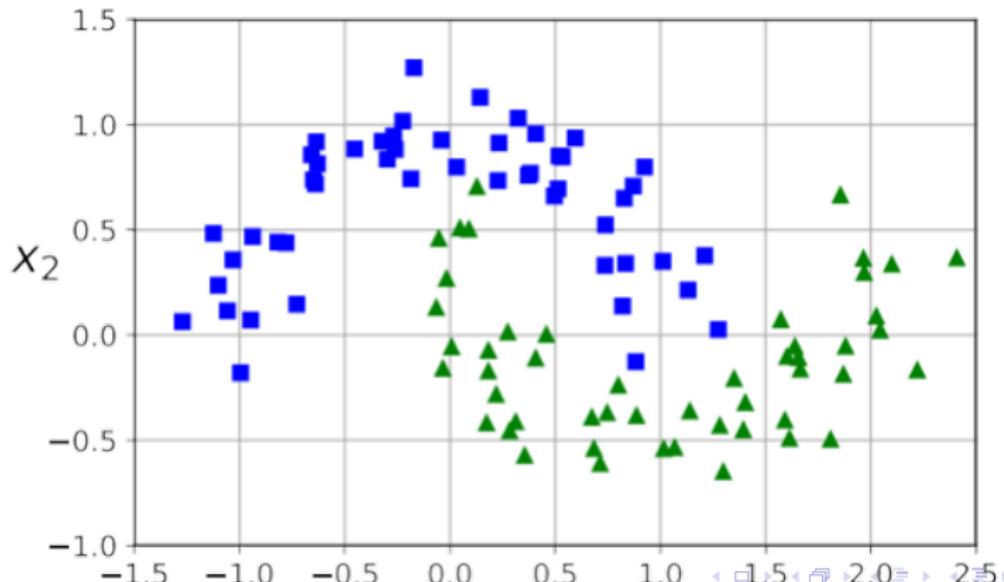


Classification Versus Clustering (1/2)

- Classification is similar to the clustering problem we saw in week 6:
 - We have a set of training data, and we wish to construct a decision function that will associate new data with sub-groups of the training data.
- The difference with classification is that the training data has already been separated into classes, and we need to generate a decision function that is consistent with these classes whenever possible, and generalises well to unseen data.
- The fact that we have class labels makes classification a supervised learning problem, whereas clustering is an unsupervised learning problem.

Classification Versus Clustering (2/2)

- The figure shows an example dataset for classification.
- It's the familiar half-moon dataset from week 6, except we now have class labels for the points in each half moon.



Iris Dataset

- Classic dataset for classification
- Classes are 3 sub-species of the iris flower
- Measurements are dimensions of different parts of the petal and sepal (labelled on the plot below)
- Good demonstration of separable and non-separable classes
- We'll use it as a running example



Iris Versicolor

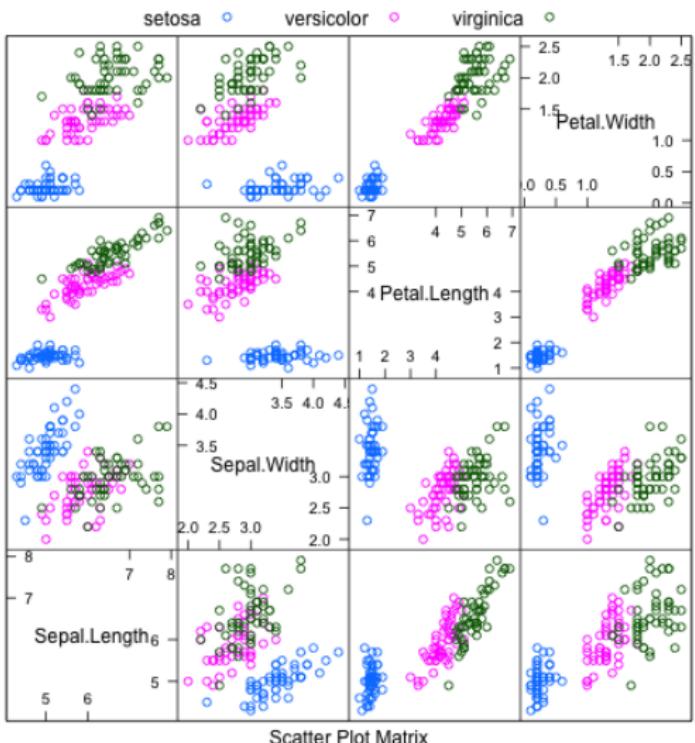


Iris Setosa



Iris Virginica

Iris Dataset Scatter Plot



Types of Classification Datasets

- Classification datasets can present different challenges...
- The simplest datasets to work with are ones in which the classes are linearly separable.
- Sometimes classes can only be separated with non-linear boundaries.
- Sometimes classes aren't separable with realistic boundaries.

Classification Dataset Examples

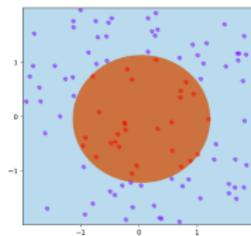
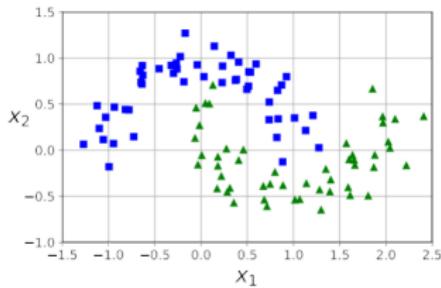
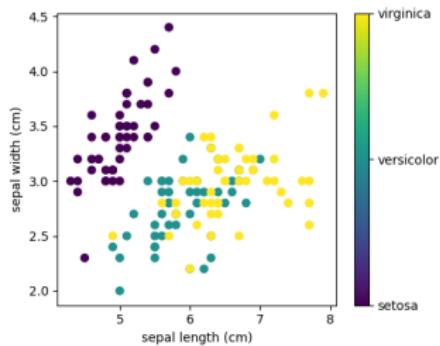
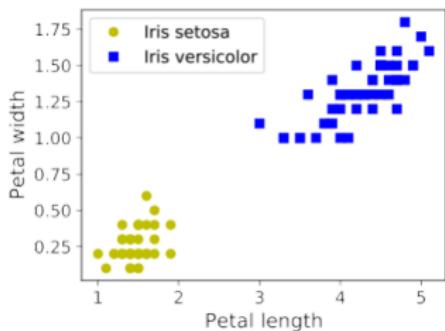


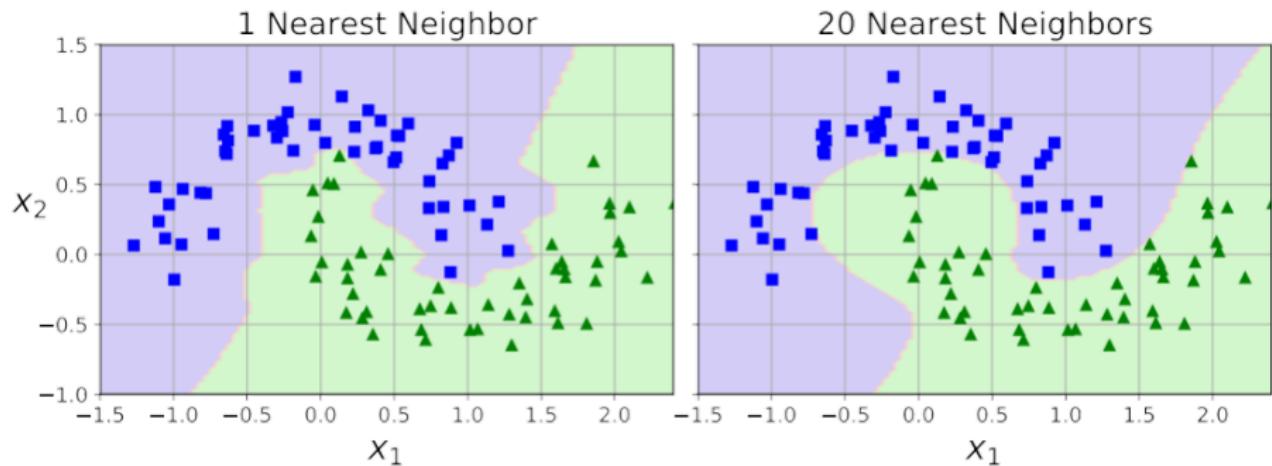
Figure: Examples of linearly separable, non-linearly separable, and inseparable

Decision Boundaries

- Classification algorithms work by constructing decision boundaries that split the sample space into regions corresponding to the different classes.
- Different classification algorithms will build different decision boundaries.
 - The plots on the next two slides show decision boundaries for the half moon dataset, created using KNN and SVM classifier (which we'll see in later slides).
- For a given classification algorithm (e.g. KNN or SVM), different choices of the algorithm's parameters will give rise to different decision boundaries (as seen in the plots).
 - Choosing the best values for these parameters is called model selection. It's out of scope for us, but see Hands-on Machine Learning if you're interested.

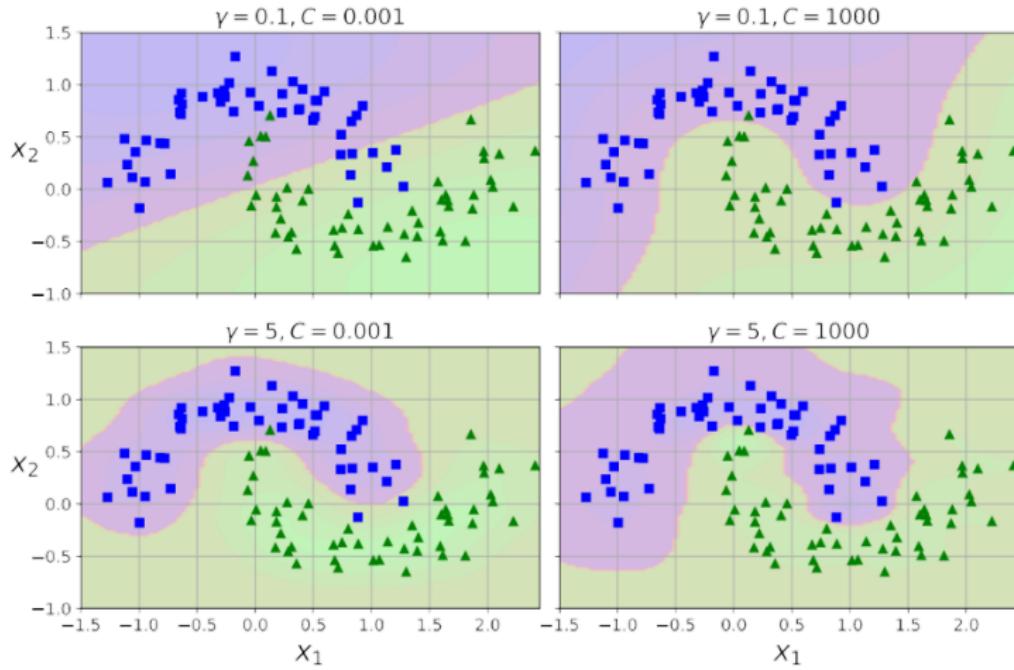
KNN Decision Boundaries

- Decision boundaries created with KNN classifiers.



SVM Decision Boundaries

- Decision boundaries created with SVM classifiers.



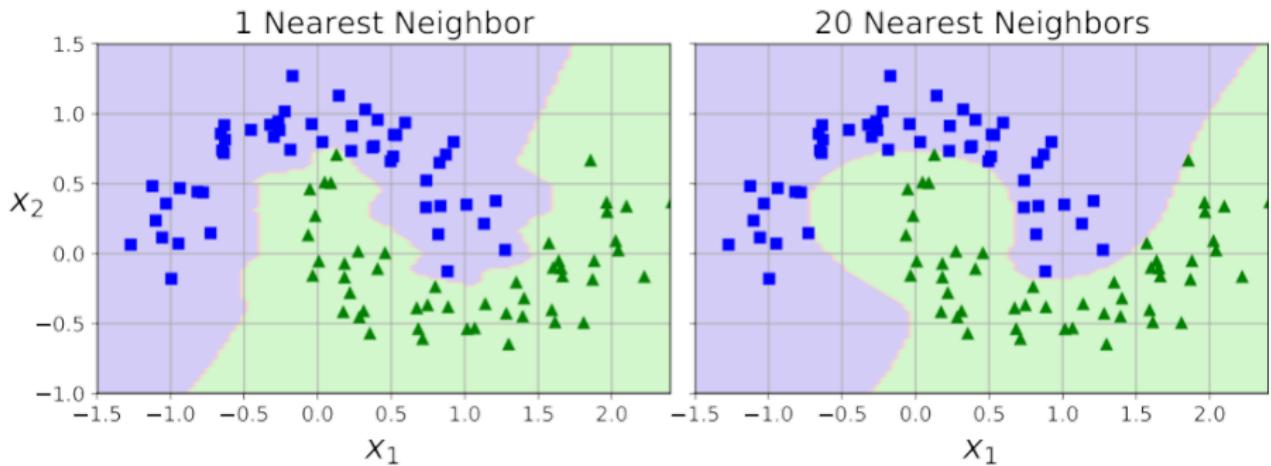
Two-Class Versus Multi-Class

- Two-class classification problems are easier to solve than multi-class, because only one decision boundary is needed.
- We'll mostly focus on two-class classifiers today, but we'll briefly introduce multi-class classifiers later in the lecture.

K-Nearest Neighbours

- **K-nearest neighbours (KNN)** is simple, intuitive decision classifier
- It assigns new points to the most common class among their k-nearest neighbours, for some k
- This often works very well for low dimensional data (when the width of the design matrix is ≤ 3)
 - **KNN is said to be the best classifier for about 1/3 of real-world classification problems** (Hastie and Tibshirani discuss this in their online ISLR lectures)
- Larger values of k give smoother decision boundaries, as shown in the plot
- **Exercise:** What's the computational complexity for this classifier?
- K-nearest neighbours unfortunately doesn't work well for high dimensional data (when the width of the design matrix is > 3)
 - The reason is outside the scope of this course. If you're interested, you can look into the "curse of dimensionality" (see either of the recommended textbooks).

KNN Decision Boundaries

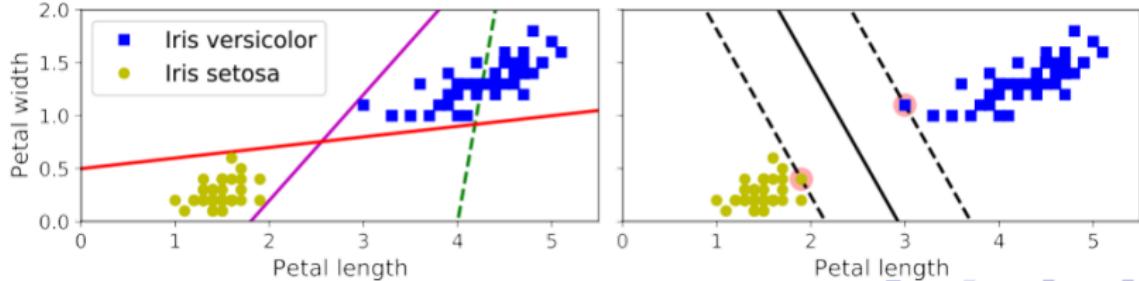


SVMs and Linear SVCs

- K-nearest neighbours doesn't work well for high dimensional data (when the width of the design matrix is large)
 - For these datasets, we need a more sophisticated approach.
- Two popular classification methods for high dimensional data are linear support vector classifiers (**SVCs**) and (nonlinear) support vector machines (**SVMs**)
- **Linear SVC** supports linear decision boundaries, and has time complexity of (almost) $O(n)$
 - Stochastic gradient descent can be used, allowing linear **SVCs** to be applied to very large datasets
- **SVM** supports non-linear decision boundaries and has time complexity of $O(n^2)$ to $O(n^3)$, depending on the specific implementation
- We'll focus on these methods for the rest of the lecture

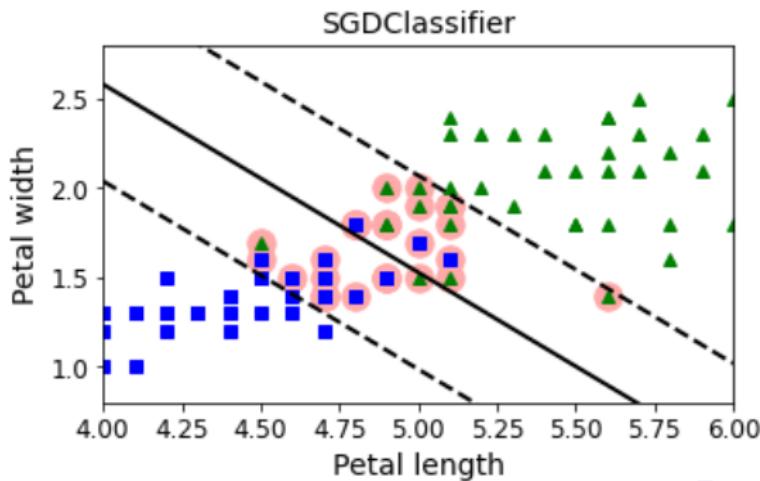
Classifiers for Linearly Separable Data

- For separable classes, it's easy to build a classifier: just draw a straight line between the two classes and treat this as the decision boundary.
- But which straight line should we choose?
 - There are infinitely many of them.
 - Two are shown in the plot below.
- One solution is to choose the separating line that permits the largest margin surrounding it which is free of observations (e.g., see the right hand plot).
- The classifier that uses this line is called the *maximal-margin classifier*, and it's used by linear **SVCs** and **SVMs**.



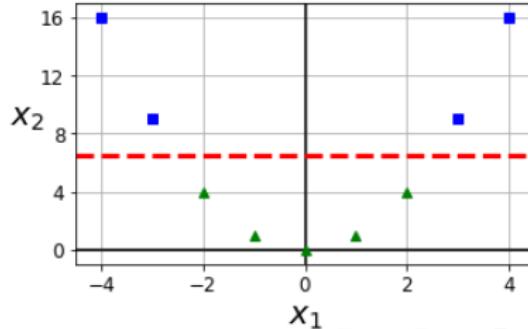
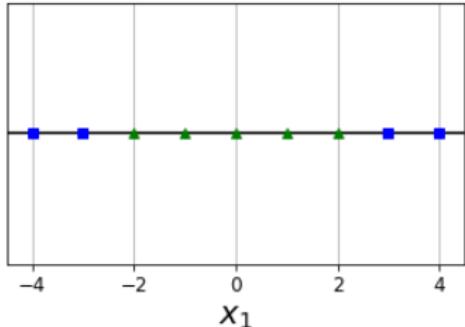
Classifiers for Linearly Separable Data

- What if the classes aren't linearly separable?
- In this case, there's no (straight) line that separates the classes.
- However, we can still use the maximal-margin classifier, if we allow a small number of observations to cross the margin.
- This classifier is called a *linear support vector classifier*.



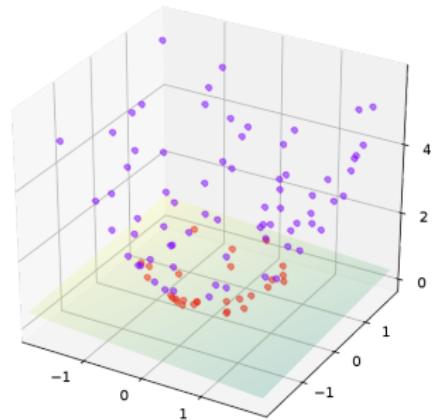
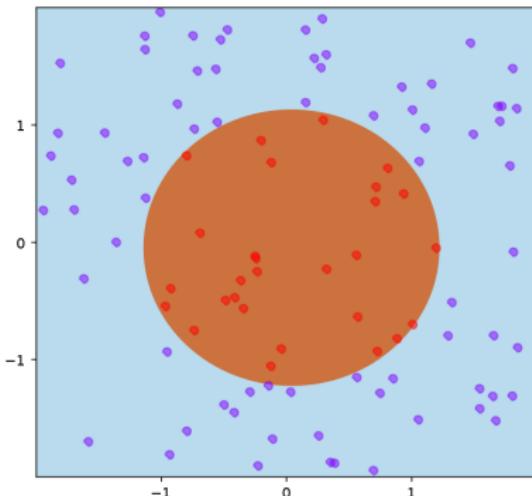
Basis Expansion

- Some classification problems just can't be handled with linear decision boundaries.
- Basis expansion is one way to deal with these datasets:
 - Augment the data with extra variables.
 - Below we turned a 1D dataset into a 2D dataset, where the second value for a given observation is set to $x_2 = x_1^2$.
- In this example, the enriched data can now be separated with a straight line in the 2D space.



Basis Expansion II

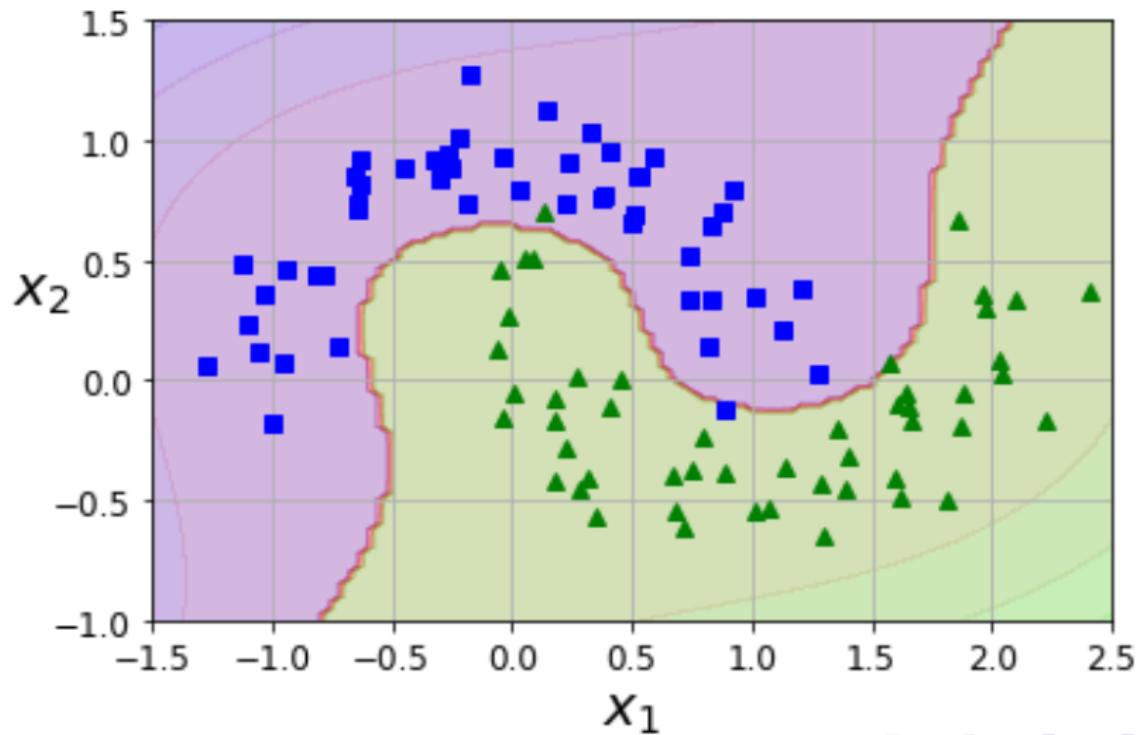
- Below is another example in which basis expansion can create a linearly separable dataset.
- This time, we added a third dimension to a 2D dataset.
- The third value for a given observation is set to $x_3 = x_1^2 + x_2^2$.
- In this case, the 3D dataset can be separated with a linear plane.



Non-Linear Decision Boundaries and SVMs

- We've seen that some classification problems just can't be handled with linear decision boundaries.
- Basis expansion works well for some, but not all, datasets.
- Enter the support vector machine...
- **SVMs** extend the linear support vector classifier to handle non-linear decision boundaries in a very general way.
- They use a technique called *the kernel trick* (which we'll see shortly):
 - This is a generalisation of the basis expansion approach that we saw in previous slides, and it can fit very flexible decision boundaries (e.g., curved or circular).
- We'll see more about the kernel trick after we cover the maths for linear support vector classifiers.

Non-Linear Decision Boundaries and SVMs



Multi-Class Classification

- Linear **SVCs** and **SVMs** are binary classifiers. However, they can be used to build multi-class classifiers as follows:
- **One-versus-the-rest (OvR)** (also called one-versus-all)
 - For an m-class dataset, fit m binary classifiers to the whole dataset, with each classifier trying to identify a single class.
 - To classify new data, run all m binary classifiers, and choose the class that achieved the highest score on its classifier.
- **One-versus-one (OvO)**
 - For an m-class dataset, fit $\frac{m(m-1)}{2}$ binary classifiers. For each pair of classes (i, j), remove all data that isn't either class i or class j. For this reduced dataset, fit a binary classifier.
 - To classify new data, see which class wins the most 1-versus-1 contests.

5

Classification

● Maths for Classifiers

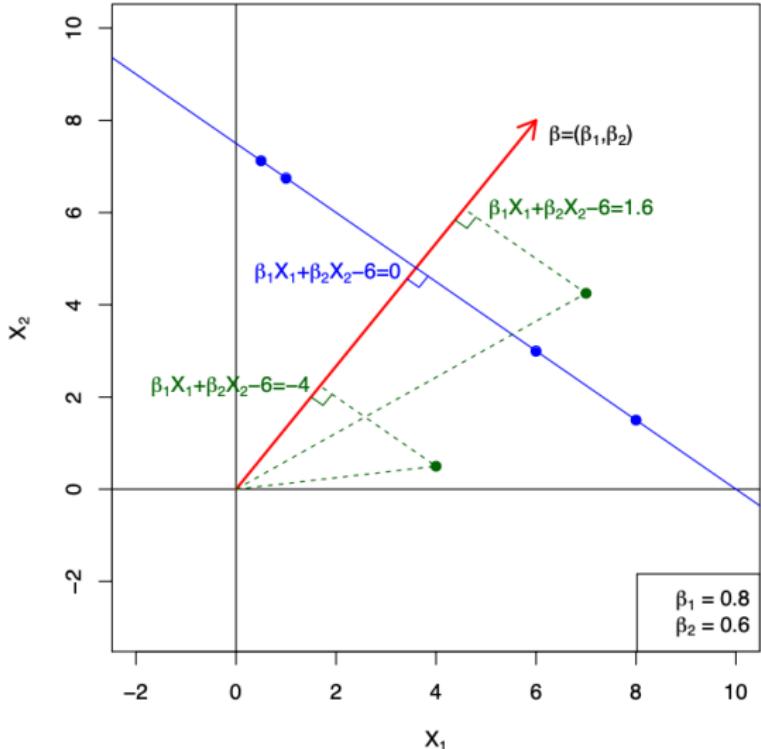
Attendance: 500240

Hyperplanes

- So far, we've seen that 2D spaces can be separated with a straight line and 3D spaces can be separated with a plane.
- In general, an N-dimensional space can be separated with an (N-1) dimensional *hyperplane*, which generalises the concept of a straight line in 2D space and a plane in 3D space.

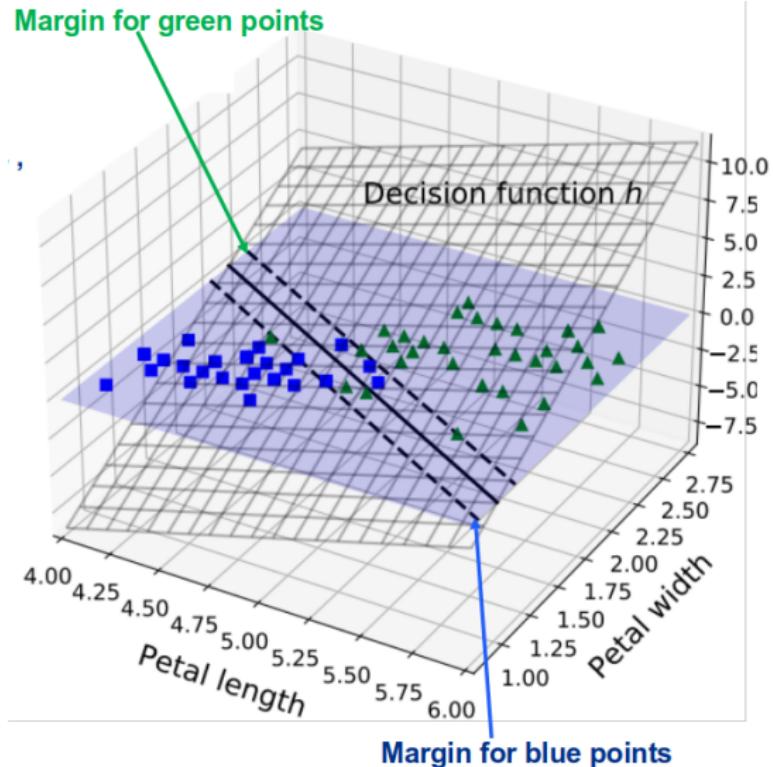
Hyperplanes

- A hyperplane is defined by a **vector and an offset**, as shown in the plot:
- We have a vector $\beta = \begin{pmatrix} \beta_1 \\ \beta_2 \end{pmatrix} = \begin{pmatrix} 0.8 \\ 0.6 \end{pmatrix}$ and an offset $\beta_0 = 6$.
- We define a hyperplane as the set of all vectors x for which $x^T \beta = \beta_0$.
- In this case, $\beta^T \beta = 1$, meaning that $x^T \beta$ is the projection of x onto the red line. (Exercise: use $x^T \beta = |x| |\beta| \cos \theta$)



Decision Functions

- For binary classification with class labels ± 1 , we can use the value of $h(\mathbf{x}) = \mathbf{x}^T \boldsymbol{\beta} - \beta_0$ as a *decision function*.
 - We decide on the class based on the sign of $h(\mathbf{x})$.
 - i.e., $\text{class}(\mathbf{x}) = \text{sign}(h(\mathbf{x}))$.
- The plot shows this decision function for iris data:
 - (The dotted lines show the **SVC** boundaries.)



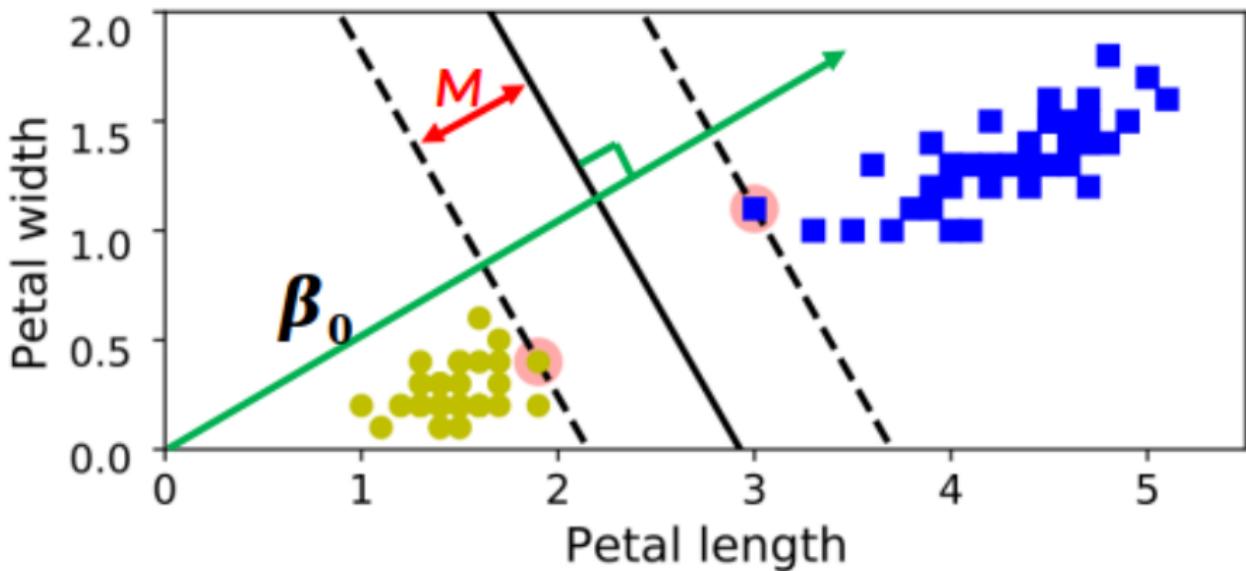
Maximal Margin Classifier

- In the slide “*Classifiers for Linearly Separable Data*”, we saw the maximal-margin classifier.
- We can find the maximal-margin classifier by solving an optimisation problem:
 - Optimisation problems were introduced in week 5.
- The maximal margin classifier is specified by a unit vector β (i.e., $\beta^T \beta = 1$) and offset β_0 , which solve the following optimisation problem:

$$\begin{aligned} & \text{maximise} && M \\ & \text{subject to} && \beta^T \beta = 1 \\ & && y_i(\beta^T \mathbf{x}_i - \beta_0) \geq M \\ & && \text{for all } i = 1, \dots, N \end{aligned}$$

- This can be rephrased as a *convex quadratic program* and solved efficiently (outside the scope of this course, but it's discussed in *Elements of Statistical Learning*).

Maximal Margin Classifier



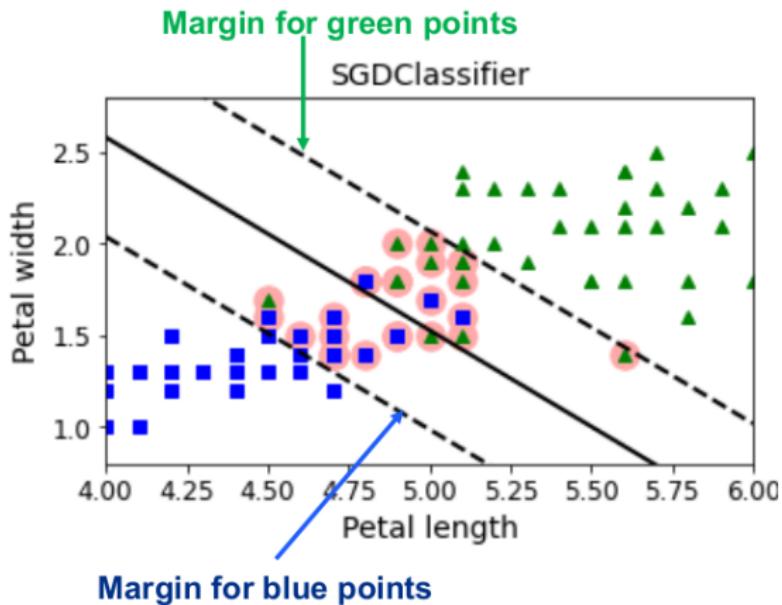
Linear Support Vector Classifier

- The Linear Support Vector Classifier requires a **small change** to the optimisation problem that we just saw.
- We allow some points to cross the separating margin (shown in **orange** in the plot on the next slide) but constrain the total distance that all points collectively pass their margins to a fixed value C :

$$\begin{aligned} & \text{maximise} && M \\ & \text{subject to} && \beta^T \beta = 1 \\ & && y_i(\beta^T \mathbf{x}_i - \beta_0) \geq M(1 - \epsilon_i) \\ & && \epsilon_i \geq 0, \quad \sum_{i=1}^N \epsilon_i \leq C \end{aligned}$$

- Again, this problem can be rephrased as a convex quadratic program and solved efficiently (outside the scope of this course, but it's discussed in *Elements of Statistical Learning*).

Linear Support Vector Classifier



Support Vectors

- It turns out that the optimal β for a linear SVC is a linear combination of the observations:

$$\beta = \sum_{i=1}^N \alpha_i \mathbf{x}_i$$

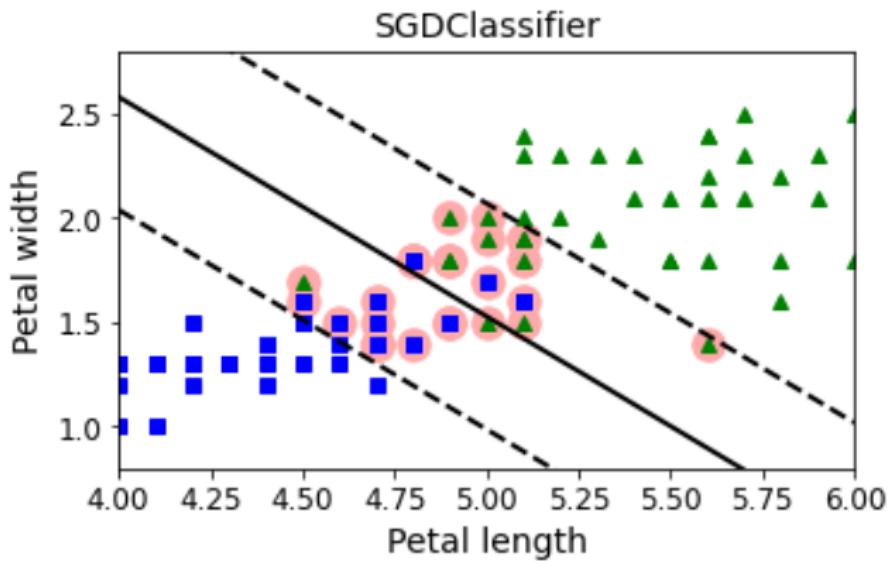
- Hence, the decision function can be written as:

$$h(\mathbf{x}) = \beta^T \mathbf{x} - \beta_0 = \left(\sum_{i=1}^N \alpha_i \mathbf{x}_i^T \mathbf{x} \right) - \beta_0$$

- Furthermore, the α_i turn out to be non-zero only for \mathbf{x}_i that are either:
 - Touching or
 - On the wrong side of their respective margins.
- These \mathbf{x}_i are referred to as the *support vectors*, shown in orange in the plot.
- The proof of this is outside the scope of the course (and quite tricky). It's given in *Elements of Statistical Learning*, which is a more advanced book from the same authors as *An Introduction to Statistical Learning*.



Support Vectors



The Kernel Trick

- We saw above that the decision function for linear support vector classifiers can be written as:

$$h_{\text{SVC}}(x) = \left(\sum_{i=1}^N \alpha_i x_i^T x \right) - \beta_0$$

where the α_i are non-zero only for x_i that are support vectors.

- We can make a small modification to this decision function, by replacing the expression $x_i^T x$ by a different function called a kernel:

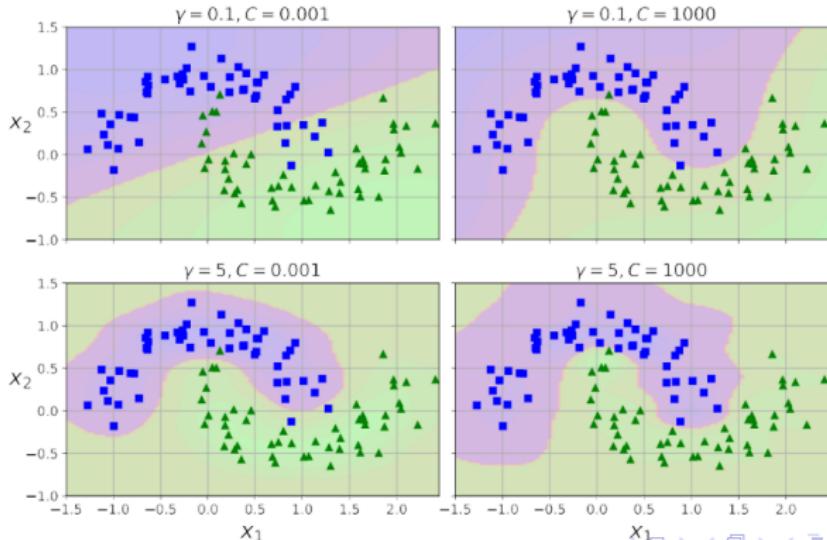
$$h_{\text{SVM}}(x) = \left(\sum_{i=1}^N \alpha_i K(x, x_i) \right) - \beta_0$$

- This small modification to the linear SVC decision function allows non-linear decision boundaries.
 - We call the resulting classifier a *Support Vector Machine*.

SVM Kernels

A popular choice of kernel function $K(x, x_i)$ is the *radial basis function* (see plot):

$$K(x, x_i) = \exp\left(-\frac{(x - x_i)^T(x - x_i)}{2\sigma^2}\right)$$



Complexity of Linear SVCs and SVMs

- Calculating the computational complexity for linear SVCs and SVMs is outside the scope of this course (because the algorithms for solving the associated optimisation problems are outside scope).
- **Time complexity:**
 - Training complexity of nonlinear SVM is generally between $O(n^2)$ and $O(n^3)$ with n the number of training instances.
 - For the linear case, the algorithm used in LinearSVC by the liblinear implementation is much more efficient than its libsvm-based SVC counterpart and can scale almost linearly to millions of samples and/or features. (Liblinear and libSVM are popular machine learning libraries)
- **Space complexity** depends on the implementation. Naive implementations have $O(n^2)$ complexity, because the entire kernel matrix is stored.
- Computational efficiency of the linear SVC can be improved with stochastic gradient descent (next slide).

Stochastic Gradient Descent

- Stochastic Gradient Descent (SGD) is a very efficient approach to fitting linear SVCs.
 - Can also be used to speed up linear regression from week 6 (see the `sklearn` class `SGDRegressor`).
- The `sklearn` class `SGDClassifier` implements a stochastic gradient descent for linear SVCs.
- SGD time complexity is $O(npI)$, where:
 - n is the number of samples (i.e., height of the design matrix),
 - p is the sample dimensionality (i.e., width of the design matrix),
 - I is the number of iterations required for convergence (which is usually small).

Machine Learning for Text Analysis

6

Machine Learning for Text Analysis

- Vector Space Document Representations
- PCA for Document Corpora
- Measuring Document Similarity
- Document Classification
- Clustering Documents

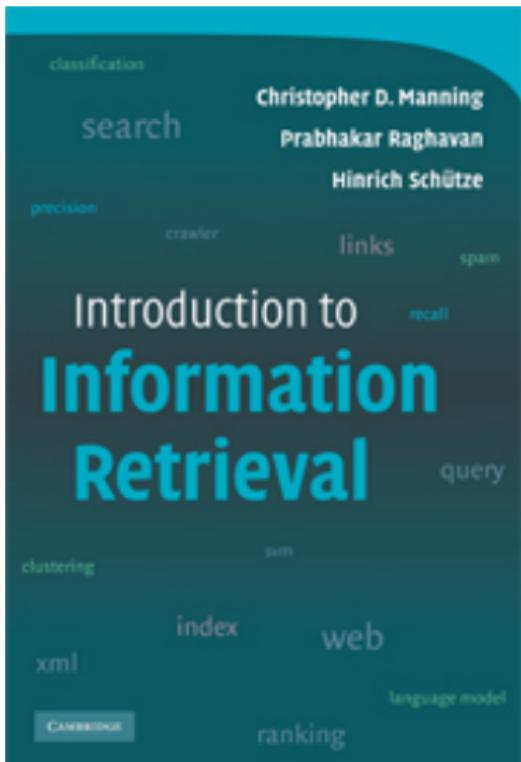
Attendance: 500240

Lecture Outline

- In this lecture, we'll look at applications of machine learning in text analysis.
- We'll apply **PCA** (week 5), **clustering** (week 6), and **classification** (week 8).
- Along the way, we'll see how to build a few different search engines.
- In order to apply machine learning for text analysis, we first need to represent text documents numerically.
- For this, we'll look at *vector space representations* of text documents. In particular:
 - **Bag of words**
 - **Term log-frequency**
 - **Term frequency, inverse document frequency (TF-IDF)**
- We'll measure the similarity of documents using their **TF-IDF** representations with a distance measure called **cosine distance**, as well as the usual Euclidean distance.

Recommended Reading

- This lecture is self-contained. However, the following might be useful:
 - **Introduction to Information Retrieval**, Manning et al.
 - **Free online**
 - Also see Christopher Manning's YouTube lectures on NLP
 - In particular, lectures 105 – 111



6

Machine Learning for Text Analysis

- **Vector Space Document Representations**
- PCA for Document Corpora
- Measuring Document Similarity
- Document Classification
- Clustering Documents

Attendance: 500240

The Data

- Training data for text analysis is typically a **corpus** (i.e., a collection, plural *corpora*) of text documents

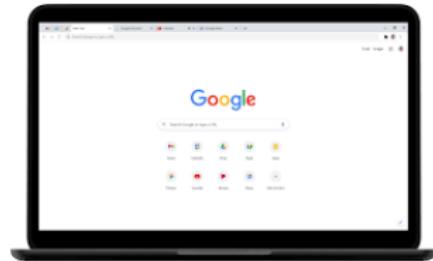
Emails, Messages, Social Media



Movie Reviews



News Articles



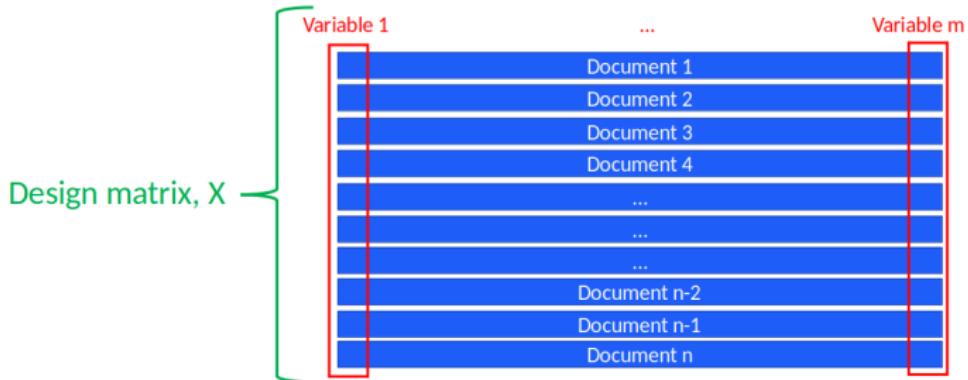
Webpages

The Goals

- How might we use these corpora?
- **Unsupervised learning:**
 - Principal component analysis
 - Topic discovery (AKA latent semantic analysis, latent semantic indexing)
- **Clustering**
 - Search enhancement
 - Keyword-free search
 - Faster search
 - Enriching search results
- **Supervised Learning:**
 - Classification
 - Sentiment analysis
 - Spam detection
 - Standing queries

Numerical Representations of Documents

- So, the training data for NLP is typically a **corpus** (i.e., a collection) of documents
- In order to use them for machine learning, we need to build a design matrix, X (week 1)
- For this, each document must be represented as a **vector** of m real variables (for some m)
- We refer to these representations as **vector space representations** of documents



Bag of Words

- One idea is to represent each document as a "bag of words".
 - For each document, form a vector of the counts for each of the (e.g.) 20,000 most commonly occurring words in English.
 - The count of occurrences of a given word t (t stands for *term*) in a document d is referred to as the **raw term frequency**, $tf_{t,d}$.
- Bag of words is referred to as a **sparse vector space representation** because most of the elements in the vector will typically be zero.



Bag of Words

- This might seem like an overly-simplified representation:
 - We discard the *ordering* of words in a document.
 - E.g., “John is quicker than Mary” has the same representation as “Mary is quicker than John”.
- In practice, bag of words representations still retain a lot of usable information, despite throwing away so much structure.
- Hence, they’re used extensively in natural language processing.
- The design matrix is:

$$X = \begin{pmatrix} tf_{0,0} & \dots & tf_{20000,0} \\ \vdots & \ddots & \vdots \\ tf_{0,n} & \dots & tf_{20000,n} \end{pmatrix}$$

Sparse and Dense Representations

- Bag of words is referred to as a **sparse vector space representation** because most of the elements in the vector will typically be zero.
- All of the representations we'll see today are sparse.
- Dense representations are outside the scope of this course, but they're complementary to sparse representations, and are often generated by neural networks such as Large Language Models (e.g. ChatGPT)

Zipf's Law and Log-Frequencies

- The bag of words representation has a drawback:
 - Common words tend to occur exponentially more often in document collections than less common words.
- This phenomenon is known as **Zipf's law**:
 - If t_1 is the most common term in the collection, t_2 is the next most common, and so on, then the collection frequency cf_i of the i th most common term is proportional to $\frac{1}{i}$.
- Hence, very common words such as “the” and “a” will tend to dominate bag of words representations.

Zipf's Law and Log-Frequencies

Hence, it's common to apply a **log-frequency weighting** to the raw term frequencies:

$$w_{t,d} = \begin{cases} 1 + \log_{10}(tf_{t,d}) & \text{if } tf_{t,d} > 0 \\ 0 & \text{otherwise} \end{cases}$$

- The “otherwise” part is necessary because $\log_{10}(0) = -\infty$.
- The design matrix is:

$$X = \begin{pmatrix} w_{0,0} & \dots & w_{20000,0} \\ \vdots & \ddots & \vdots \\ w_{0,n} & \dots & w_{20000,n} \end{pmatrix}$$



Inverse Document Frequency

- Log-frequencies are a step in the right direction, but we can do better still...
- Rare terms tend to be more informative than frequent terms, so we'd like to make them more prominent in the document representation.
- We want to give high weights for rare terms and lower rates for frequent terms.
- We can measure the commonality of a word with **document frequency** df_t :
 - The document frequency of a term is the number of documents that contain the word at least once.

Inverse Document Frequency

- In practice, we use the so-called **inverse document frequency**:

$$idf_t = \log \left(\frac{N}{df_t} \right)$$

- N is the number of documents in the corpus.
- $\frac{df_t}{N}$ is the fraction of documents in which term t appears.
- Log is used to dampen the effect of idf_t .
- Using $\frac{N}{df_t}$ itself would be too strong a factor.

TF-IDF

- Combining the log frequency weighting with the inverse document frequency weighting, we get a value called the *term frequency, inverse document frequency* (TF-IDF):

$$w_{t,d} = \begin{cases} (1 + \log tf_{t,d}) \times \log \left(\frac{N}{df_t} \right) & \text{if } tf_{t,d} > 0 \\ 0 & \text{otherwise} \end{cases}$$

- Again, the design matrix is:

$$X = \begin{pmatrix} w_{0,0} & \dots & w_{20000,0} \\ \vdots & \ddots & \vdots \\ w_{0,n} & \dots & w_{20000,n} \end{pmatrix}$$

- This is the most commonly-used sparse vector space representation for documents.



6

Machine Learning for Text Analysis

- Vector Space Document Representations
- **PCA for Document Corpora**
- Measuring Document Similarity
- Document Classification
- Clustering Documents

Attendance: 500240

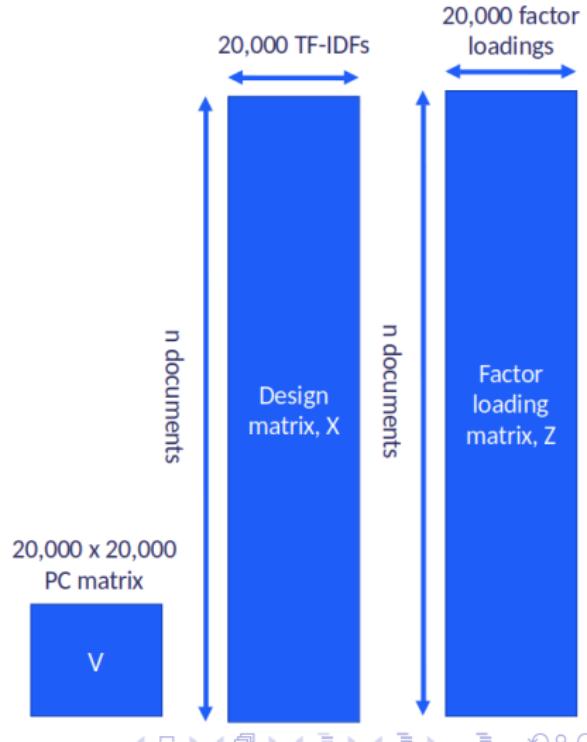
Latent Semantic Analysis

- We can apply PCA to “term log-frequency” design matrices.
 - In the jargon of natural language processing, this is called **latent semantic analysis (LSA)**, or sometimes **latent semantic indexing (LSI)**.
- When we do this, the principal components have an interesting interpretation...
 - Recall from weeks 5 and 6 that applying PCA to images of faces gave principal components that represented things like glasses, smiles, cheekbones, face width, nose length.
 - Applying PCA to term log-frequency design matrices gives principal components which represent document topics.

Latent Semantic Analysis

- Recall that the principal component matrix is $\mathbf{V} = (\mathbf{v}_1 \quad \mathbf{v}_2 \quad \dots \quad \mathbf{v}_{20,000})$, where \mathbf{v}_i is the i^{th} principal component.
- As usual with PCA, \mathbf{X} , \mathbf{V} and \mathbf{Z} are related by the equations:

$$\mathbf{Z} = \mathbf{X}\mathbf{V}, \quad \mathbf{X} = \mathbf{Z}\mathbf{V}^T$$



6

Machine Learning for Text Analysis

- Vector Space Document Representations
- PCA for Document Corpora
- **Measuring Document Similarity**
- Document Classification
- Clustering Documents

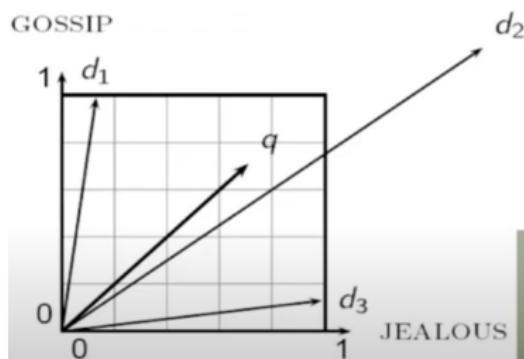
Attendance: 500240

Euclidean and Cosine Distances

- So far, we've seen a few vector space representations of documents, and we saw that the TF-IDF vector space representation is a particularly sensible one.
- How should we measure the distance between documents using their TF-IDF representations?
- E.g., imagine that we have a document q that reads “jealous and gossip”, and we'd like to compare it to three other documents:
 - d_1 is a short document mostly about gossip.
 - d_2 is a long document about both jealousy and gossip.
 - d_3 is a short document mostly about jealousy.

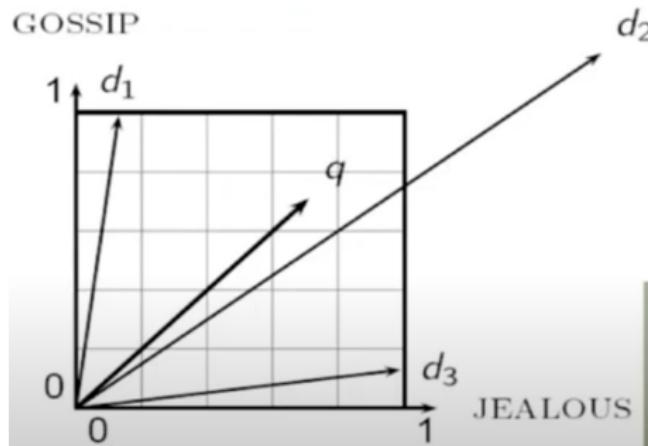
Euclidean and Cosine Distances

- q is closer in Euclidean distance to d_1 and d_3 than d_2 .
- So, Euclidean distance isn't a great measure of similarity in this case.
- In this instance, a better measure of similarity would be the angle between vectors.
- We call this **cosine similarity**.
- **Exercise:** Consider the Euclidean and cosine distances between any document d and the document formed by appending d to itself. Which would be a more sensible measure of distance in this case?



Search Engines

- Cosine distance between TF-IDF document representations is one way to build a search engine.
- E.g., imagine that q from the previous example were a search term and d_1 , d_2 , and d_3 were webpages.



6

Machine Learning for Text Analysis

- Vector Space Document Representations
- PCA for Document Corpora
- Measuring Document Similarity
- **Document Classification**
- Clustering Documents

Attendance: 500240

Classification and the Contiguity Hypothesis

- In the last lecture, we saw a few algorithms for classification:
 - K-nearest neighbours
 - Linear support vector classifiers
 - Support vector machines
- For classification of documents, one should use a document representation that approximately obeys the **contiguity hypothesis**:
 - Documents in the same class form a contiguous region in the representation space, and regions of different classes do not overlap.
 - TF-IDF generally provides good contiguity.
 - See Manning for more details if you're interested.
- Both SVMs and linear SVCs use Euclidean distance for classification, so for these algorithms, one would use Euclidean distance for TF-IDF vectors.
- KNN could support either Euclidean distance or cosine distance between TF-IDF vectors.

Sentiment Analysis on Movie Reviews

- In order to train a classifier, we would need a body of documents which already has class labels.
- One standard example is a collection of movie reviews:
 - E.g., the Rotten Tomatoes movie review dataset on Kaggle.



Sentiment Analysis on Movie Reviews

Classify the sentiment of sentences from the Rotten Tomatoes dataset



Kaggle · 860 teams · 7 years ago

6

Machine Learning for Text Analysis

- Vector Space Document Representations
- PCA for Document Corpora
- Measuring Document Similarity
- Document Classification
- Clustering Documents

Attendance: 500240

Clustering Documents

- We saw in week 6 that clustering can be applied to any dataset for which we have a dissimilarity measure.
 - We also saw that k-means only supports one dissimilarity measure: Euclidean distance.
- Today we have seen 2 ways to measure the dissimilarity between documents:
 - Euclidean distance between TF-IDF vectors.
 - Cosine distance between TF-IDF vectors.
- Hence, we can use these dissimilarity measures to cluster text documents with the usual clustering algorithms such as k-means, HAC, and DBSCAN.

Cluster Hypothesis for Information Retrieval

- So, we can cluster text documents. But why would we want to?
- Clustering can be useful for document search.
- This application is motivated by a truism known as the **Cluster Hypothesis for Information Retrieval**:
 - *“Documents in the same cluster behave similarly with respect to information needs”.*
 - E.g., documents in the same cluster are relevant to the same search query.
 - See Manning for more discussion if you’re interested.

Search Result Clustering

- Web search engines often retrieve more webpages than a user can realistically view.
- Clustering can be applied to these webpages to help the user.
- E.g., the Vivisimo search engine:

The screenshot shows the Vivisimo search interface. At the top, there's a navigation bar with links for company, products, solutions, demos, partners, and press. Below that is a search bar containing the query "data mining". To the right of the search bar is a dropdown menu set to "Search the Web" and a blue "Search" button. A yellow banner below the search bar displays the message "Query data mining returned 221 documents." On the left, a sidebar titled "Clustered Results" lists categories: data mining (221), Knowledge Discovery (36), Analysis (29), Machine learning (18), Conference, Pacific-Asia (3), Research (3), Weka 3 (2), Qub (2), DataMining, Knowledge Discovery And Machine Learning (2), Other Topics (5), and Papers (16). The main content area shows two sponsored links: "Find Data Mining Solutions" and "Logfile Analysis With Metaposition". Both links have a "Sponsored Link" label to their right. Below these are several organic search results, starting with a link to "KDnuggets: Data Mining, Web Mining, and Knowledge Discovery Guide".

company | products | solutions | demos | partners | press

data mining Search the Web Search

Query data mining returned 221 documents.

[Find Data Mining Solutions](#) [New Window] [Preview] Sponsored Link
Free online research tool. Search 22,000 enterprise solution profiles to find the right Data Mining solution fast. Registration Required. - <http://www.knowledgestorm.com/Keywords/lookup?x=Data%20Minin...>

[Logfile Analysis With Metaposition](#) [New Window] [Preview] Sponsored Link
Combine logfile and ranking analysis with Metaposition and gain valuable information on new traffic chances. - <http://www.metaposition.com/index.php?logfile>

1. [KDnuggets: Data Mining, Web Mining, and Knowledge Discovery Guide](#) [New Window]
... KnowledgeDiscovery, Genomic Mining , Web Mining Data Mining Consulting | Data Mining Jobs | Advertising | Site Map ... KDnuggets News, the leading newsletter on Data Mining and Knowledge... More results from: www.kdnuggets.com May 6, 2003- 9 KB 2. [Data Mining and Analytic Technologies](#) (Kurt Thearling)... kurt@thearl... Data Mining , if you haven't heard of it before ... creation of a number of commercial data mining software applications, including ... an award winning system of data mining ...
URL: <http://www.kdnuggets.com/>
Sources: Lycos 1, Looksmart 8, Netscape 1, MSN 1, Altavista 1

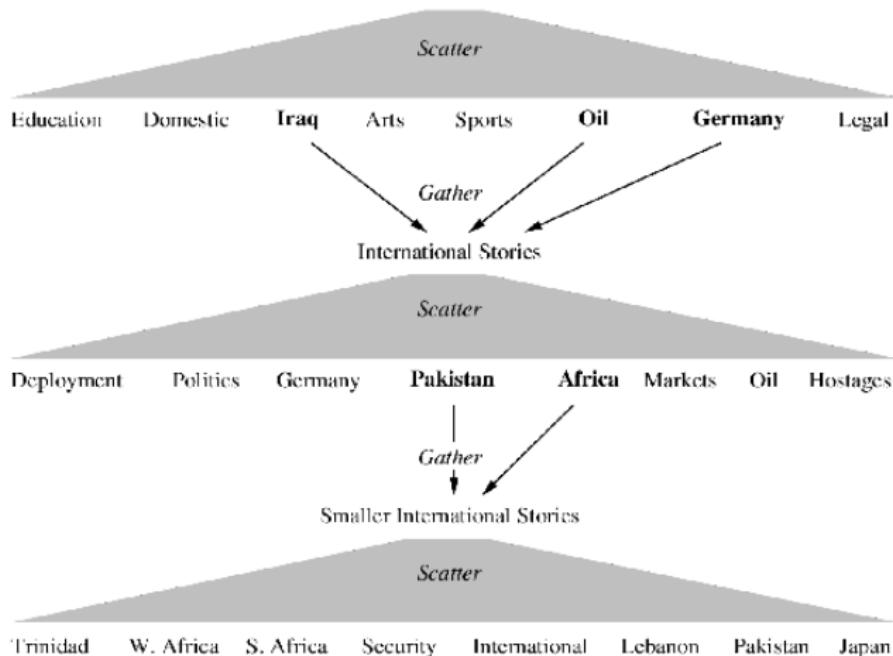
Dom Richards (University of Liverpool) November 11, 2024 235 / 239

Big Data Analytics

Scatter-Gather

- **Scatter-Gather** uses clustering to provide an alternative to keyword search for navigating large bodies of text.
- The scatter-gather algorithm begins by splitting the document corpus into clusters.
- The user then picks a subset of the clusters.
- The user-selected clusters are gathered together and clustered again.
- This process is repeated until a cluster of interest is found.

Scatter-Gather



Collection Clustering

- Clustering can be used to enrich the results produced by a traditional search engine.
- E.g., imagine you have created a search engine that selects documents based on the cosine similarity between the TF-IDF representations of the user's search term and the documents in your corpus.
- For documents that were selected by your TF-IDF search engine, documents from the same cluster can also be added to the search results.
- This will enrich the search result by adding documents which are related to the search query but don't contain exactly the same words.

Cluster-Based Retrieval

- Clustering can be used to speed up search.
- Rather than searching the whole document corpus, find the clusters that most closely match the query and only consider the documents in those clusters.
- Offers a good speedup for a small decrease in the expected search quality.