

Lecture Notes on
Modern Data Science

Module 06: Cloud Computing - Spark

Gang Li
School of Information Technology
Deakin University, VIC 3125, Australia



TULIP Tool for Universal Learning and Intelligent Processing

SIT742 | Modern Data Science (G. Li)

1

Unit Learning Outcomes

- ULO1:**
 - Develop knowledge of and discuss new and emerging fields in data science.
 - Cloud Computing and Hadoop/Spark Platform
- ULO4:**
 - Use appropriate platform to collect and process relatively large datasets.
 - Programming in Spark

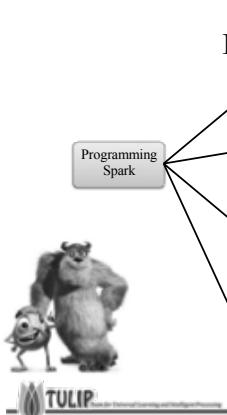


TULIP Tool for Universal Learning and Intelligent Processing

SIT742 | Modern Data Science (G. Li)

1

Road map



```

graph LR
    PS[Programming Spark] --> HO[Hadoop Overview]
    PS --> CC[Cloud Computing]
    PS --> SRA[Spark Runtime Architecture]
    PS --> C[Closures]
    HO --> MR[MapReduce]
    CC --> MR
    CC --> RDD[RDDs]
    CC --> PRD[Pair RDD]
    SRA --> BV[Broadcast Variable]
    SRA --> ACU[Accumulator]
    C --> BV
    C --> ACU
  
```

TULIP Tool for Universal Learning and Intelligent Processing

SIT742 | Modern Data Science (G. Li)

2

Hadoop Essentials



- Hadoop EcoSystem
- HDFS
- MapReduce



TULIP Tool for Universal Learning and Intelligent Processing

SIT742 | Modern Data Science (G. Li)

3

Data Intensive Computing

- Data intensive computing
 - Concerns with the production, manipulation and analysis of data in the range of hundreds of megabytes (MB) to petabytes (PB) and beyond
 - A range of supporting parallel and distributed computing technologies to deal with the challenges of data representation, reliable shared storage, efficient algorithms and scalable infrastructure to perform analysis



TULIP Tool for Universal Learning and Intelligent Processing

SIT742 | Modern Data Science (G. Li)

4

Data Intensive Computing

10 GB
100 GB
1 TB



TULIP Tool for Universal Learning and Intelligent Processing

SIT742 | Modern Data Science (G. Li)

5

Data Intensive Computing

10 TB

100 TB

XX PB

TULIP Data for Universal Learning and Intelligent Processing

SIT742 | Modern Data Science (G. Li)

6

Data Intensive Computing

- Problems with Scale-Up
 - Disk streaming speed ~ 50MB/s
 - 3TB ~ 17.5 Hours
 - 1PB ~ 8 months
 - Seeking speeds even worse than r/w
 - ~10ms for a seek
 - Enough time to read 0.5MB

TULIP Data for Universal Learning and Intelligent Processing

SIT742 | Modern Data Science (G. Li)

7

Meet the Hadoop

- What is Hadoop?
 - It's a framework for running applications on large clusters of commodity hardware which produces huge data and to process it
 - Based on 2004 Google MapReduce Paper

TULIP Data for Universal Learning and Intelligent Processing

SIT742 | Modern Data Science (G. Li)

8

Meet the Hadoop

- HPC attitude
 - “The problem of disk-limited, loosely-coupled data analysis was solved by throwing more disks and using weak scaling”
- Flip-side:
 - A single novice developer can write real, scalable, 1000+ node data-processing tasks in Hadoop-family tools in an afternoon
 - MPI... less so

TULIP Data for Universal Learning and Intelligent Processing

SIT742 | Modern Data Science (G. Li)

9

Meet the Hadoop

- Hadoop 2+
 - “An increasingly important tool for quantitative researchers to have at their disposal”
 - Spark, Pregel/Giraph for general, MPI-like communication patterns
 - Enormous data processing capabilities, batch + interactive

TULIP Data for Universal Learning and Intelligent Processing

SIT742 | Modern Data Science (G. Li)

10

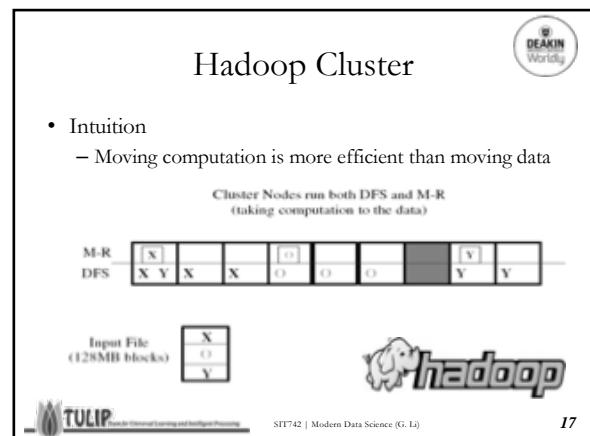
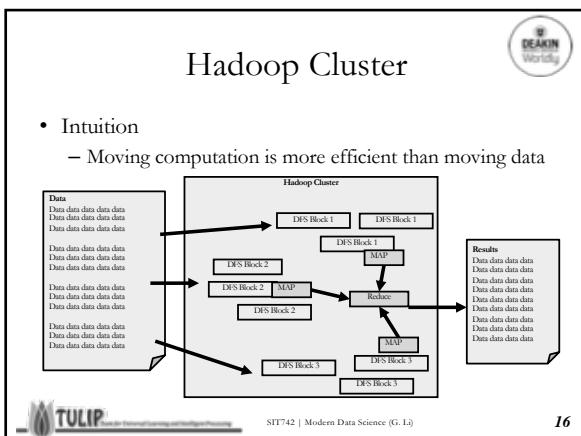
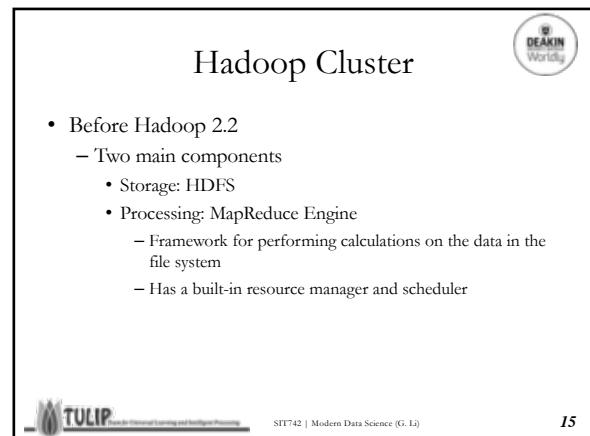
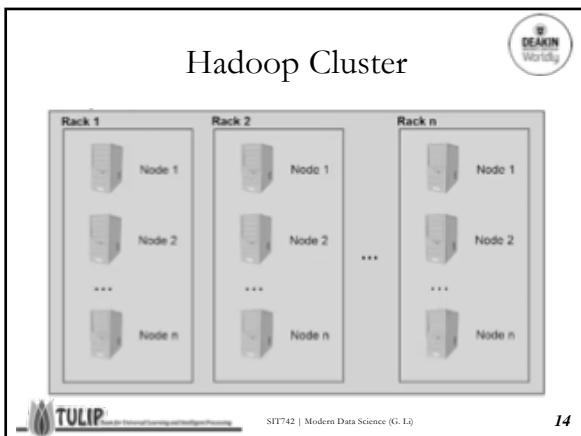
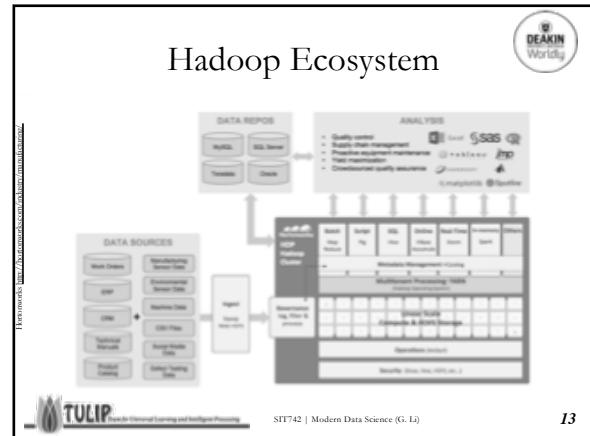
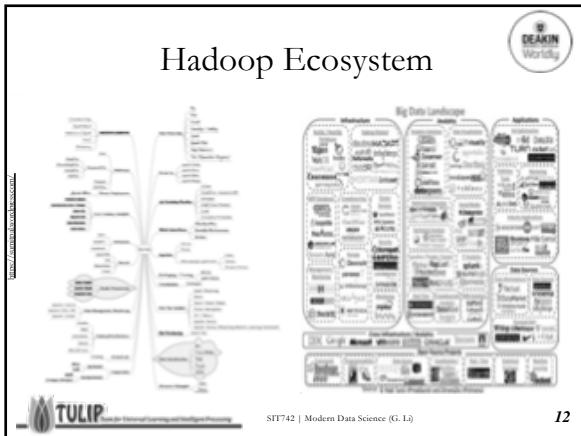
Everything is Converging

- Good ideas are good ideas
- These models are all converging at the largest scales
 - MPI is trying to grow fault tolerance
 - RDBM are scaling up
 - Hadoop is growing tightly coupled
 - Much faster than MPI codes are growing fault tolerance
 - Emerging tools tightly coupled computation atop Hadoop

TULIP Data for Universal Learning and Intelligent Processing

SIT742 | Modern Data Science (G. Li)

11



Hadoop Distributed File System (HDFS)

- A distributed file system designed to run on commodity hardware
 - Share similarities with existing distributed file systems and supports traditional hierarchical file organization
 - Reliable data replication and accessible via Web interface and Shell commands
 - Benefits:* Fault tolerant, high throughput, streaming data access, robustness and handling of large data sets
 - HDFS is not a general purpose F/S



SIT742 | Modern Data Science (G. Li)

18



Worldly



Hadoop Distributed File System (HDFS)

- What HDFS is not good for?
 - Low-latency reads
 - High-throughput rather than low latency for small chunks of data
 - Large amount of small files
 - Better for millions of large files instead of billions of small files
 - Multiple writers
 - Single write per file
 - Writes only at the end of file, no support for arbitrary offset



SIT742 | Modern Data Science (G. Li)

19

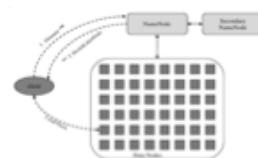
Hadoop Distributed File System (HDFS)

- DFS Master “**Namenode**”
 - Manages the filesystem namespace
 - Maintain file map to list blocks + location mappings
 - Manages block allocation/replication
 - Checkpoints namespace and journals namespace changes for reliability
 - Control access to namespace
- DFS Slaves “**Datanodes**” handle block storage
 - Stores blocks using the underlying OS’s files
 - Clients access the blocks directly from datanodes
 - Periodically sends block reports to Namenode and check block integrity
- DFS “**Secondary Namenode**”
 - for checkpointing instead of backup



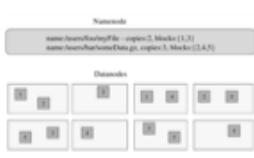
SIT742 | Modern Data Science (G. Li)

20



Hadoop Distributed File System (HDFS)

- Files are split into blocks
- Blocks
 - Single unit of storage: a contiguous piece of information on a disk
 - Transparent to user
 - Managed by Namenode, stored by Datanode
 - Blocks are either 64MB or 128MB



SIT742 | Modern Data Science (G. Li)

21

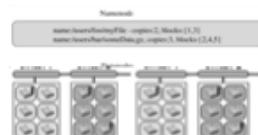
Hadoop Distributed File System (HDFS)

- Same block is replicated on multiple machines (3 by default)
 - Replicate placement are rack aware
 - 1st replica on the local rack
 - 2nd replica on the local rack but different machine
 - 3rd replica on the different rack
- Namenode determines replica placement



SIT742 | Modern Data Science (G. Li)

22



Hadoop Distributed File System (HDFS)

- DataNode
 - Each block replica on a DataNode is represented by two files in the local native filesystem.
 - The first file contains the data itself
 - The second file records the block’s metadata including checksums for the data and the generation stamp.
 - At startup each DataNode connects to a NameNode and performs a handshake.
 - The handshake verifies that the DataNode is part of the NameNode and runs the same version of software
 - A DataNode identifies block replicas in its possession to the NameNode by sending a block report.



SIT742 | Modern Data Science (G. Li)

23

Hadoop Distributed File System (HDFS)

- NameNode
 - The NameNode does not directly send requests to DataNodes. It uses replies to heartbeats to send instructions to the DataNodes
 - The instructions include commands to replicate blocks to other nodes, remove local block replicas, re-register and send an immediate block report, and shut down the node



SIT742 | Modern Data Science (G. Li)

24

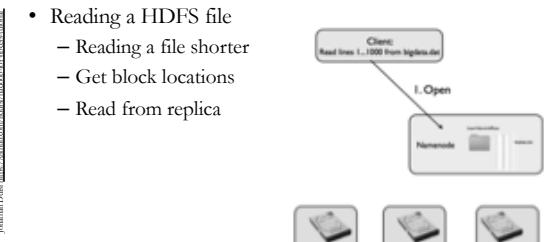


Worldly

Hadoop Distributed File System (HDFS)

- Reading a HDFS file
 - Reading a file shorter
 - Get block locations
 - Read from replica

JOURNAL DRAFT - DEAKIN UNIVERSITY



SIT742 | Modern Data Science (G. Li)

25

Hadoop Distributed File System (HDFS)

- Reading a HDFS file
 - Reading a file shorter
 - Get block locations
 - Read from replica



SIT742 | Modern Data Science (G. Li)

26

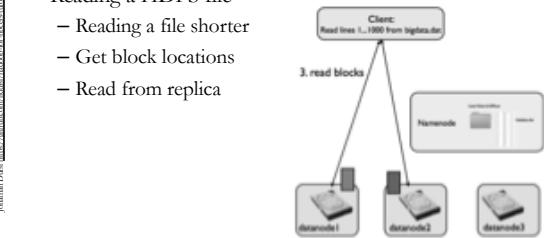


Worldly

Hadoop Distributed File System (HDFS)

- Reading a HDFS file
 - Reading a file shorter
 - Get block locations
 - Read from replica

JOURNAL DRAFT - DEAKIN UNIVERSITY

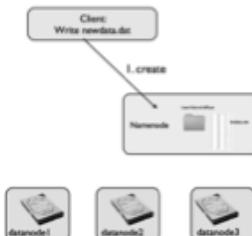


SIT742 | Modern Data Science (G. Li)

27

Hadoop Distributed File System (HDFS)

- Writing a HDFS file
 - Creating file
 - Get nodes for blocks
 - Start writing
 - Data nodes coordinate replication
 - Get Acknowledgement
 - Complete



SIT742 | Modern Data Science (G. Li)

28

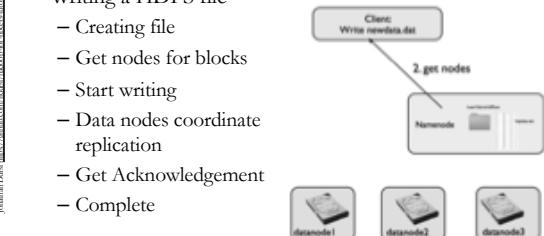


Worldly

Hadoop Distributed File System (HDFS)

- Writing a HDFS file
 - Creating file
 - Get nodes for blocks
 - Start writing
 - Data nodes coordinate replication
 - Get Acknowledgement
 - Complete

JOURNAL DRAFT - DEAKIN UNIVERSITY

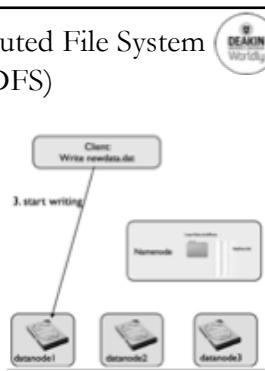


SIT742 | Modern Data Science (G. Li)

29

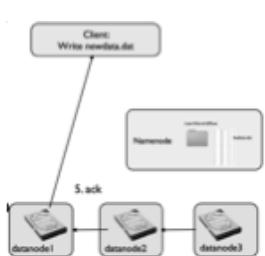
Hadoop Distributed File System (HDFS)

- Writing a HDFS file
 - Creating file
 - Get nodes for blocks
 - Start writing
 - Data nodes coordinate replication
 - Get Acknowledgement
 - Complete



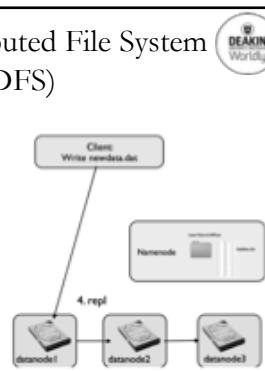
Hadoop Distributed File System (HDFS)

- Writing a HDFS file
 - Creating file
 - Get nodes for blocks
 - Start writing
 - Data nodes coordinate replication
 - Get Acknowledgement
 - Complete



Hadoop Distributed File System (HDFS)

- Writing a HDFS file
 - Creating file
 - Get nodes for blocks
 - Start writing
 - Data nodes coordinate replication
 - Get Acknowledgement
 - Complete



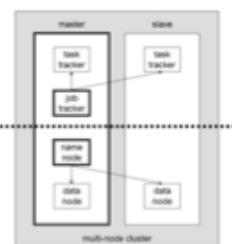
Hadoop Server Roles



SIT742 | Modern Data Science (G. Li)

34

Hadoop Server Roles

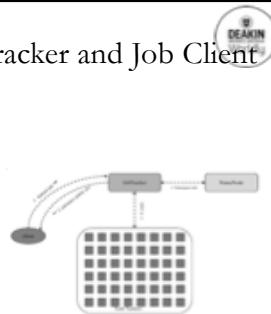


SIT742 | Modern Data Science (G. Li)

35

Job Tracker, Task Tracker and Job Client

- Map/Reduce Master “Jobtracker”
 - Accepts MR jobs submitted by users
 - Assigns Map and Reduce tasks to Tasktrackers
 - Monitors task and tasktracker status, re-executes tasks upon failure
- Map/Reduce Slaves “Tasktrackers”
 - Run Map and Reduce tasks upon instruction from the Jobtracker
 - Manage storage and transmission of intermediate output.



SIT742 | Modern Data Science (G. Li)

36

MapReduce Paradigm

- Programming model developed at Google
 - Sort/merge based distributed computing
 - The underlying system takes care of the partitioning of the input data, scheduling the program's execution across several machines, handling machine failures, and managing required inter-machine communication.
 - This is the key for Hadoop's success



SIT742 | Modern Data Science (G. Li)

37

MapReduce Paradigm

- Mapping workers to Processors
 - The input data (on HDFS) is stored on the local disks of the machines in the cluster.
 - The MapReduce master takes the location information of the input files into account and attempts to schedule a map task on a machine that contains a replica of the corresponding input data.
 - Failing that, it attempts to schedule a map task near a replica of that task's input data.
 - When running large MapReduce operations on a significant fraction of the workers in a cluster, most input data is read locally and consumes no network bandwidth.



SIT742 | Modern Data Science (G. Li)

38

MapReduce Paradigm

- Mapping workers to Processors
 - The input data (on HDFS) is stored on the local disks of the machines in the cluster.
 - The MapReduce master takes the location information of the input files into account and attempts to schedule a map task on a machine that contains a replica of the corresponding input data.
 - Failing that, it attempts to schedule a map task near a replica of that task's input data.
 - When running large MapReduce operations on a significant fraction of the workers in a cluster, most input data is read locally and consumes no network bandwidth.



SIT742 | Modern Data Science (G. Li)

39

MapReduce Paradigm

- Mapping workers to Processors
 - The input data (on HDFS) is stored on the local disks of the machines in the cluster.
 - The MapReduce master takes the location information of the input files into account and attempts to schedule a map task on a machine that contains a replica of the corresponding input data.
 - Failing that, it attempts to schedule a map task near a replica of that task's input data.
 - When running large MapReduce operations on a significant fraction of the workers in a cluster, most input data is read locally and consumes no network bandwidth.



SIT742 | Modern Data Science (G. Li)

40

MapReduce Paradigm

- Mapping workers to Processors
 - The input data (on HDFS) is stored on the local disks of the machines in the cluster.
 - The MapReduce master takes the location information of the input files into account and attempts to schedule a map task on a machine that contains a replica of the corresponding input data.
 - Failing that, it attempts to schedule a map task near a replica of that task's input data.
 - When running large MapReduce operations on a significant fraction of the workers in a cluster, most input data is read locally and consumes no network bandwidth.



SIT742 | Modern Data Science (G. Li)

41

MapReduce Paradigm

- Task Granularity
 - The map phase has M pieces and the reduce phase has R pieces.
 - M and R should be much larger than the number of **worker machines**.
 - Having each worker perform many different tasks improves dynamic load balancing, and also speeds up recovery when a worker fails.
 - Larger the M and R, more the decisions the master must make
 - R is often constrained by users because the output of each reduce task ends up in a separate output file.
 - Typically, (at Google), M = 200,000 and R = 5,000, using 2,000 worker machines.



SIT742 | Modern Data Science (G. Li)

42

Hadoop 2.2+

- YARN
 - MapReduce V2
 - DataNodes still exist
 - JobTrackers and TaskTrackers no longer exist



SIT742 | Modern Data Science (G. Li)

43

Hadoop 2.2+

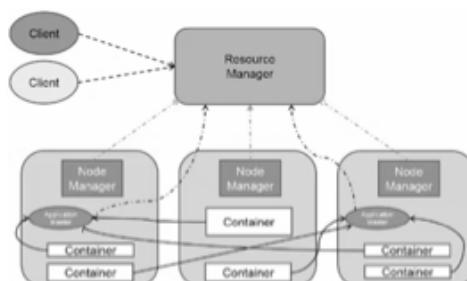
- YARN
 - Provides generic scheduling and resource management
 - More than just batch processing
 - More efficient scheduling and workload management
 - No more balancing between map and reduce slots



SIT742 | Modern Data Science (G. Li)

44

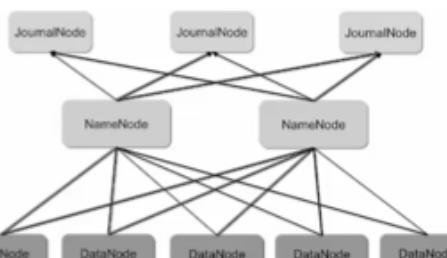
Hadoop High Availability



SIT742 | Modern Data Science (G. Li)

45

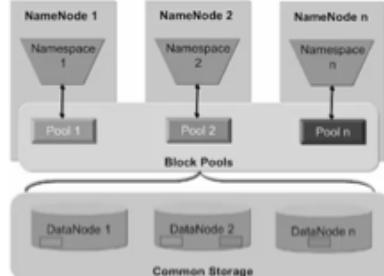
Hadoop High Availability



SIT742 | Modern Data Science (G. Li)

46

Hadoop High Availability



SIT742 | Modern Data Science (G. Li)

47

Map Reduce



- Map + Reduce
- Apache Spark
- History Review

TULIP Tutorials for Universal Learning and Intelligent Processing

SIT742 | Modern Data Science (G. Li)

48

MapReduce

- Main Operations
 - Map
 - Reduce
- Input:
 - Key Value pairs
 - Write once and Read Many (WORM) data
- Output:
 - Key Value pairs
 - Written to disk



TULIP Tutorials for Universal Learning and Intelligent Processing

SIT742 | Modern Data Science (G. Li)

49

MapReduce

- An Example Task:
 - Counts lines in all *.txt files that match <regex> and displays the counts in descending order

File 1 File 2 Result

TULIP Tutorials for Universal Learning and Intelligent Processing

SIT742 | Modern Data Science (G. Li)

50

MapReduce

Map tasks

(0, C) → [(C, 1)]	(A, [1]) → (A, 1)
(2, B) → []	(C, [1, 1, 1]) → (C, 3)
(4, B) → []	
(6, C) → [(C, 1)]	
(0, C) → [(C, 1)]	
(2, A) → [(A, 1)]	

Reduce tasks

File 1 File 2 Result

TULIP Tutorials for Universal Learning and Intelligent Processing

SIT742 | Modern Data Science (G. Li)

51

MapReduce

Map tasks

(0, C) → [(C, 1)]	(A, [1]) → (A, 1)
(2, B) → []	(C, [1, 1, 1]) → (C, 3)
(4, B) → []	
(6, C) → [(C, 1)]	
(0, C) → [(C, 1)]	
(2, A) → [(A, 1)]	

Reduce tasks

>>> grep -E 'A|C' *.txt | sort | uniq -c | sort -nr

TULIP Tutorials for Universal Learning and Intelligent Processing

SIT742 | Modern Data Science (G. Li)

52

MapReduce

Hadoop is an implementation of the Map-Reduce framework for distributed processing of large data sets.

Map tasks

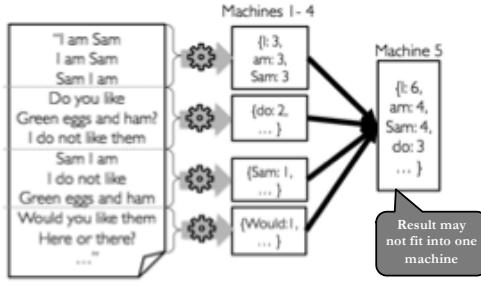
Reduce tasks

TULIP Tutorials for Universal Learning and Intelligent Processing

SIT742 | Modern Data Science (G. Li)

53

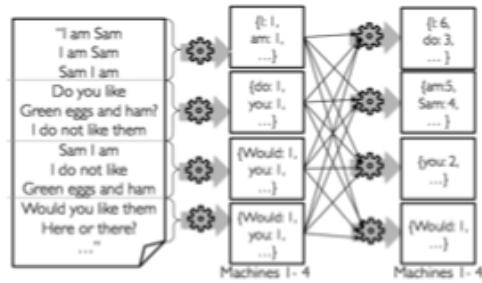
What if the Document is Really Big?



SIT742 | Modern Data Science (G. Li)

54

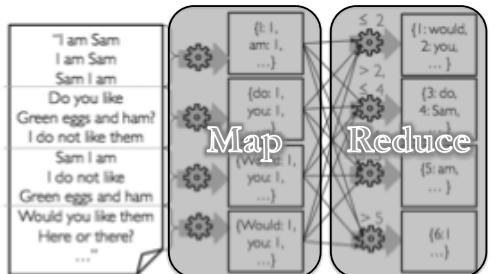
What if the Document is Really Big?



SIT742 | Modern Data Science (G. Li)

55

What if the Document is Really Big?



SIT742 | Modern Data Science (G. Li)

56

What if the Document is Really Big?

- Need to handle more data?
 - Just add more Mappers/Reducers!
- No need to handle multi-threaded code
 - Mappers and Reducers are typically single threaded and deterministic
 - Determinism allows for restarting of failed jobs
 - Mappers/Reducers run entirely independent of each other
 - In Hadoop, they run in separate JVMs

SIT742 | Modern Data Science (G. Li)

57

How to Deal with Failures?

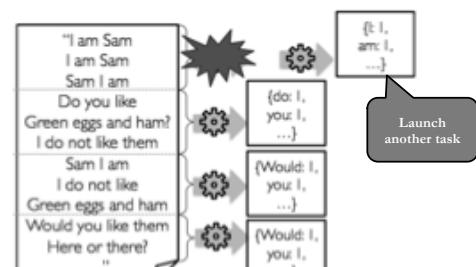
- Failures are norm in commodity hardware
 - 1 server fails every 3 years → with 10k nodes 10 faults/day
- **Worker failure**
 - Detect failure via periodic heartbeats
 - Re-execute completed and in-progress map tasks
 - Task completion committed through master
- **Master failure**
 - Single point of failure; Resume from Execution Log
- Robust
 - Google's experience:
 - lost 1600 of 1800 machines once!, but finished fine



SIT742 | Modern Data Science (G. Li)

58

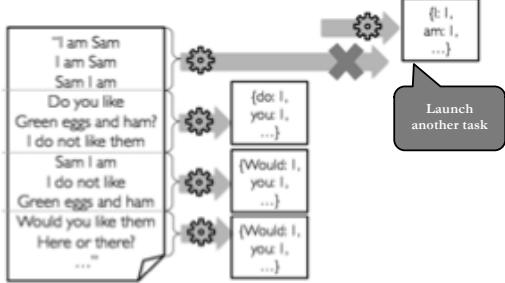
How to Deal with Failures?



SIT742 | Modern Data Science (G. Li)

59

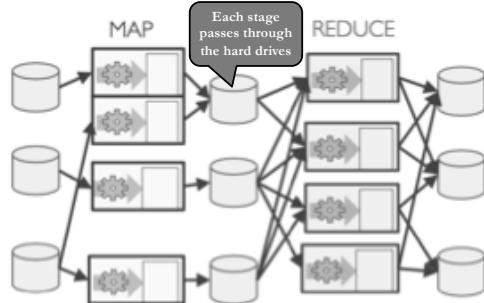
How to Deal with Slow Tasks?



SIT742 | Modern Data Science (G. Li)

60

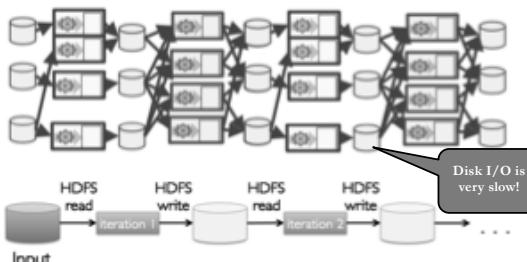
Map Reduce: Distributed Execution



SIT742 | Modern Data Science (G. Li)

61

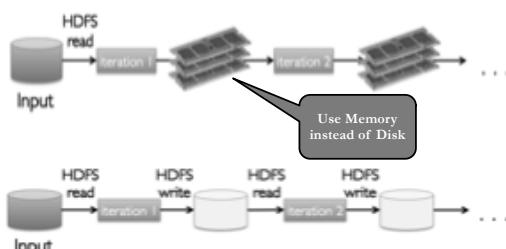
Map Reduce: Iterative Jobs



SIT742 | Modern Data Science (G. Li)

62

Map Reduce: Iterative Jobs



SIT742 | Modern Data Science (G. Li)

63

Apache Spark

- The **Spark** Computing Framework
 - Keep more data in-memory
 - A new distributed execution engine that provides programming abstraction and parallel runtime to hide complexities of fault-tolerance and slow machines



SIT742 | Modern Data Science (G. Li)

64

Apache Spark

- Resilient Distributed Dataset (RDD)**
 - a read-only, partitioned collection of records spread across a cluster, stored in memory or on disk
 - RDDs can only be created through deterministic operations on either
 - data in stable storage or
 - other RDDs.



SIT742 | Modern Data Science (G. Li)

65

Spark vs Hadoop

DEAKIN Worldwide

Interest over time

May 2011 - May 2016

TULIP Data for Universal Learning and Intelligent Processing

SIT742 | Modern Data Science (G. Li)

66

Spark vs Hadoop MapReduce

DEAKIN Worldwide

	Spark	Hadoop MapReduce
Storage	In-memory or on disk	Disk only
Operations	Map, Reduce, Join, Sample, etc.	Map and Reduce
Execution Model	Batch, Interactive, Streaming	Batch
Programming Environment	Scala, Java, R, and Python	Java

TULIP Data for Universal Learning and Intelligent Processing

SIT742 | Modern Data Science (G. Li)

67

A Record in Daytona GraySort

DEAKIN Worldwide

Hadoop MR Record	Spark Record	Spark
102.5 TB	100 TB	1 PB
Elapsed Time	72 mins	23 mins
# Nodes	2100	206
# Cores	50400 physical	6592 virtualized
Cluster disk throughput [est.]	3150 GB/s	618 GB/s
Sort Benchmark Daytona Rules	Yes	Yes
Network	dedicated data center, 10Gbps	virtualized (EC2) 10Gbps network
Sort rate	1.42 TB/min	4.27 TB/min
Sort rate/node	0.67 GB/min	20.7 GB/min
		22.5 GB/min

TULIP Data for Universal Learning and Intelligent Processing

SIT742 | Modern Data Science (G. Li)

68

History Review

DEAKIN Worldwide

GFS Paper MapReduce Paper

Spark Paper RDD Paper

2002 2004 2006 2008 2010 2012 2014 2016

MapReduce @ Google Hadoop @ Yahoo! Hadoop Summit Spark@UCB Apache Spark Summit

TULIP Data for Universal Learning and Intelligent Processing

SIT742 | Modern Data Science (G. Li)

69

Spark Essentials

DEAKIN Worldwide

- Spark Runtime Architecture
- Spark Context
- Resilient Distributed Dataset

TULIP Data for Universal Learning and Intelligent Processing

SIT742 | Modern Data Science (G. Li)

70

Spark Runtime Architecture

DEAKIN Worldwide

- In distributed mode, Spark uses a master/slave architecture
 - one central coordinator (**driver**) and many distributed workers (**executors**)
- Spark Application contains
 - The driver runs in its own Java process
 - Each executor is a separate Java process

TULIP Data for Universal Learning and Intelligent Processing

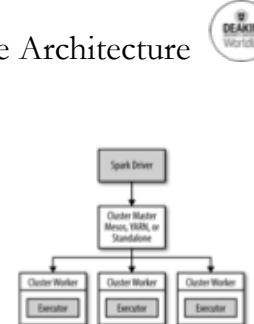
SIT742 | Modern Data Science (G. Li)

71

Spark Runtime Architecture

- The Driver**

- The process where method `main()` runs
- Two duties
 - Converting a user program into tasks
 - Send works to executors as tasks
 - Scheduling tasks on executors
 - Coordinate the scheduling of individual tasks on executors
 - Complete view of application's executors

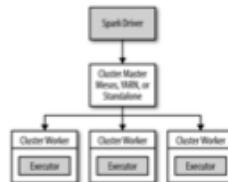


72

Spark Runtime Architecture

- Executors**

- The worker processes responsible for running individual tasks
- Two duties
 - Running the tasks and return results to the driver
 - Providing in-memory storage for RDDs that are cached, through the Block Manager.

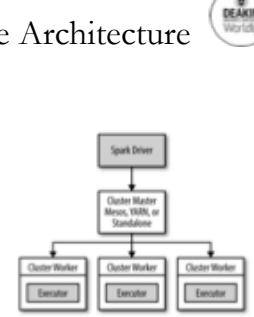


73

Spark Runtime Architecture

- Cluster Manager**

- Standalone cluster manager, or external managers such as YARN, Mesos, etc.
- Launch executors on behalf of the driver program
- Launch the driver sometimes
- Terms
 - **Master and Worker**
 - Centralized and distributed portions of the cluster manager
 - **Driver and Executor**
 - The processes that execute each Spark application

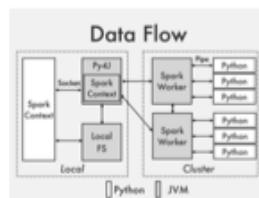


74

Python Spark (PySpark)

- PySpark** is built on top of Spark's Java API

- Data is processed in Python and cached/shuffled in JVM
- In the Python driver, `SparkContext` uses Py4J to launch a JVM and create a `JavaSparkContext`.
 - Py4J is only used on the driver for local communication between the Python and Java `SparkContext` objects;
 - large data transfers are performed through a different mechanism.



SIT742 | Modern Data Science (G. Li)

75

Spark Context

- Driver programs access Spark through a `SparkContext` object
 - which tells Spark how and where to access a cluster
- Use `SparkContext` to create RDDs



76

Initializing a SparkContext

- The master parameter for a `SparkContext` determines which type and size of cluster to use

Value	Explanation
<code>spark://host:port</code>	Connect to a Spark Standalone cluster at the specified port. By default Spark Standalone masters use port 7077.
<code>mesos://host:port</code>	Connect to a Mesos cluster master at the specified port. By default Mesos masters listen on port 5050.
<code>yarn</code>	Connect to a YARN cluster. When running on YARN you'll need to set the <code>HADOOP_CONF_DIR</code> environment variable to point the location of your hadoop configuration directory, which contains information about the cluster.
<code>local</code>	Run in local mode with a single core.
<code>local[N]</code>	Run in local mode with N cores.
<code>local[*]</code>	Run in local mode and use as many cores as the machine has.

SIT742 | Modern Data Science (G. Li)

77

Initializing a SparkContext

- The master parameter for a **SparkContext** determines which type and size of cluster to use

```
from pyspark import SparkConf, SparkContext
conf = SparkConf().setMaster("local").setAppName("SIT742")
sc = SparkContext(conf = conf)
```

– PySpark is not on **sys.path** by default, we can use

```
!pip install findspark
import findspark
findspark.init()

import pyspark
sc = pyspark.SparkContext(appName="SIT742")
```

© Deakin University 2017. All rights reserved.

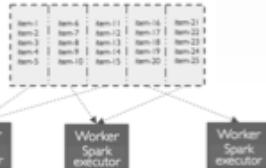


SIT742 | Modern Data Science (G. Li)

78

Resilient Distributed Datasets

- RDD is an **immutable** distributed collection of objects
 - Each RDD is split into multiple partitions, which may be computed on different nodes of the cluster
 - Track lineage to efficiently re-compute the lost data
 - Programmer specifies the number of partitions for an RDD

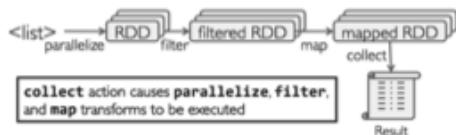


TULIP: Tutorials for Unsupervised Learning and Intelligent Processing

79

Resilient Distributed Datasets

- Working with RDDs
 - Create an RDD
 - Apply transformations to an RDD: map, filter, etc.
 - Apply actions to an RDD: collect, count



© Deakin University 2017. All rights reserved.



SIT742 | Modern Data Science (G. Li)

80

Resilient Distributed Dataset

- Creating RDDs
- RDD Operations
 - Transformation
 - Actions
- Caching RDDs



TULIP: Tutorials for Unsupervised Learning and Intelligent Processing

SIT742 | Modern Data Science (G. Li)

81

Creating RDDs

- Two ways to create RDDs:
 - Parallelizing** existing Python collections
 - We can specify the number of partitions (4 here)
 - No computation occurs with **sc.parallelize**
 - Spark only records how to create the RDD with 4 partitions here



```
>>> data = [1, 2, 3, 4, 5]
>>> rdd = sc.parallelize(data, 4)
>>> rdd
ParallelCollectionRDD[6] at parallelize at PythonRDD.scala:392
```

© Deakin University 2017. All rights reserved.



SIT742 | Modern Data Science (G. Li)

82

Creating RDDs

- Two ways to create RDDs:
 - Loading** an dataset from HDFS, or other storage
 - We can specify the number of partitions (4 here)
 - No computation occurs with **sc.parallelize**
 - Spark only records how to create the RDD with 4 partitions here

```
>>> fdd = sc.textfile("/path/to/filename", 4)
>>> fdd
MapPartitionsRDD[14] at textFile at
NativeMethodAccessImpl.java:-2
```

TULIP: Tutorials for Unsupervised Learning and Intelligent Processing

SIT742 | Modern Data Science (G. Li)

83

RDD Operations

<u>Transformations</u>	<u>Actions</u>
<ul style="list-style-type: none"> Operations on RDDs that return a new RDD Lazy Evaluation <ul style="list-style-type: none"> Results not computed right away Spark remembers set of transformations applied to base dataset 	<ul style="list-style-type: none"> Operations that return a result to the driver program or write it to storage Eager <ul style="list-style-type: none"> Cause Spark to kick off a computation Mechanism for getting results out of Spark

TULIP Data for Universal Learning and Intelligent Processing SIT742 | Modern Data Science (G. Li) 84

RDD Operations: Transformation

- Python **lambda** Functions
 - Small anonymous functions that is
 - Restricted to a single expression
 - Not bound to a name
 - Can use **lambda** functions wherever function objects are required

```
lambda x: x*x
lambda x: x!=1
lambda x: [x, x+5]
```

TULIP Data for Universal Learning and Intelligent Processing SIT742 | Modern Data Science (G. Li) 85

RDD Operations: Transformation

- Element-wise transformations
 - Transformation **filter()**
 - Takes in a function and returns an RDD that only has elements that pass the **filter()** function

```
Rdd = sc.parallelize([1,2,3,4])
Rdd.filter(lambda x: x!=1)
```

TULIP Data for Universal Learning and Intelligent Processing SIT742 | Modern Data Science (G. Li) 86

RDD Operations: Transformation

- Element-wise transformations
 - Transformation **map()**
 - Takes in a function and applies it to each element in the RDD with the result of the function being the new value of each element in the resulting RDD

```
Rdd = sc.parallelize([1,2,3,4])
Rdd.map(lambda x: x*x)
```

TULIP Data for Universal Learning and Intelligent Processing SIT742 | Modern Data Science (G. Li) 87

RDD Operations: Transformation

- Element-wise transformations
 - Transformation **flatMap()**
 - Get back an RDD that consists of the elements from all of its iterators

```
tokenizer('coffee panda') = List('coffee', 'panda')
Rdd1 = sc.parallelize(['coffee panda', 'happy panda'])
Rdd1.flatMap(tokenizer)
```

TULIP Data for Universal Learning and Intelligent Processing SIT742 | Modern Data Science (G. Li) 88

RDD Operations: Transformation

- Sampling Transformation
 - sample()** an RDD
 - We can specify with or without replacement, or the fraction

```
Rdd = sc.parallelize([1,2,3,3])
sdd = Rdd.sample(False, 0.5)
```

TULIP Data for Universal Learning and Intelligent Processing SIT742 | Modern Data Science (G. Li) 89

RDD Operations: Transformation

Table 3-2. Basic RDD transformations on an RDD containing {1, 2, 3, 3}

Function name	Purpose	Example	Result
map()	Apply a function to each element in the RDD and return an RDD of the result.	<code>rdd.map(x => x + 1)</code>	{2, 3, 4, 4}
flatMap()	Apply a function to each element in the RDD and return an RDD of the contents of the iterators returned. Often used to extract words.	<code>rdd.flatMap(x => x.toInt() * 3)</code>	{3, 3, 3, 3}
filter()	Return an RDD consisting of only elements that pass the condition passed to <code>filter(_)</code> .	<code>rdd.filter(x => x >= 1)</code>	{1, 2, 3}
distinct()	Remove duplicates.	<code>rdd.distinct()</code>	{1, 2, 3}
sample(withReplacement, fraction, seed)	Sample an RDD, with or without replacement.	<code>rdd.sample(False, 0.5)</code>	Nondeterministic



SIT742 | Modern Data Science (G. Li)

90

RDD Operations: Transformation

- Pseudo Set Operations

- RDDs support many operations of mathematical sets:

- **distinct, union, intersect, subtract**
- All expensive except **union** because they involve shuffling



SIT742 | Modern Data Science (G. Li)

91

RDD Operations: Transformation

Table 3-3. Two-RDD transformations on RDDs containing {1, 2, 3} and {3, 4, 5}

Function name	Purpose	Example	Result
union()	Produce an RDD containing elements from both RDDs.	<code>rdd.union(other)</code>	{1, 2, 3, 3, 4, 5}
intersection()	RDD containing only elements found in both RDDs.	<code>rdd.intersection(other)</code>	{3}
subtract()	Remove the contents of one RDD (e.g., remove training data).	<code>rdd.subtract(other)</code>	{1, 2}
carteresian()	Cartesian product with the other RDD.	<code>rdd.cartesian(other)</code>	{(1, 3), (1, 4), (2, 3)}



SIT742 | Modern Data Science (G. Li)

92

RDD Operations: Action

- **reduce()**

- takes a function that operates on two elements of the type in the RDD and returns a new element of the same type
- should be **commutative** and **associative** so that it can be computed correctly in parallel

```
>>> Rdd = sc.parallelize([1,2,3])
>>> Rdd.reduce(lambda a, b: a+b)
Value: 6
```



SIT742 | Modern Data Science (G. Li)

93

RDD Operations: Action

- **fold()**

- takes a function similarly as **reduce()** does, but takes a “zero value” to be used for initial call on each partition
- should be the identity element for the operation
 - 0 for +, 1 for *, etc.
- return type the same as RDD elements

```
>>> Rdd = sc.parallelize([1,2,3])
>>> Rdd.fold(0)(lambda a, b: a+b)
Value: 6
```



SIT742 | Modern Data Science (G. Li)

94

RDD Operations: Action

- **aggregate()**

- we can also supply an initial zero value of the type we want to return
- a 1st function to combine the elements from RDD with the accumulator
- a 2nd function to merge two accumulators given that each node accumulates its own results locally

```
sumCount = Rdd.aggregate((0,0),
                         (lambda acc, value: (acc[0]+value, acc[1]+1),
                          (lambda acc1, acc2: (acc1[0]+acc2[0], acc1[1]+acc2[1]))))
return sumCount[0]/float(sumCount[1])
```



SIT742 | Modern Data Science (G. Li)

95

RDD Operations: Action

- collect()**
 - return the entire RDD's contents to the driver program
- take()**
 - returns n elements from the RDD
 - attempts to minimize the number of partitions it accesses, so may be biased

```
>>> Rdd = sc.parallelize([1,2,3])
>>> Rdd.take(2)
Value: [1,2]
>>> Rdd.collect()
Value: [1,2,3]
```



SIT742 | Modern Data Science (G. Li)

96

RDD Operations: Action

- top(n)**
 - return the top n elements of the RDD
- count()**
 - returns the number of elements in the RDD

```
>>> Rdd = sc.parallelize([1,2,3])
>>> Rdd.count()
Value: 3
```



SIT742 | Modern Data Science (G. Li)

97

RDD Operations: Action

Table 3-4. Basic actions on an RDD containing [1, 2, 3, 3]

Function name	Purpose	Example	Result
collect()	Return all elements from the RDD.	rdd.collect()	(1, 2, 3, 3)
count()	Number of elements in the RDD.	rdd.count()	4
countByValue()	Number of times each element occurs in the RDD.	rdd.countByValue()	(1, 1), (2, 1), (3, 2)
take(num)	Return num elements from the RDD.	rdd.take(2)	(1, 2)
top(num)	Return the top num elements from the RDD.	rdd.top(2)	(3, 3)
takeOrdered(num)(ordering)	Return num elements based on provided ordering.	rdd.takeOrdered(2)(myOrdering)	(3, 3)



SIT742 | Modern Data Science (G. Li)

98

RDD Operations: Action

takeSample(withReplacement, num, [seed])	Return num elements at random.	rdd.takeSample(false, 1)	Non-deterministic
reduce(func)	Combine the elements of the RDD together in parallel (e.g., sum).	rdd.reduce((x, y) == x + y)	9
fold(zeroValue)(func)	Same as reduce() but with the provided zero value.	rdd.fold(0)((x, y) => x + y)	9
aggregate(zeroValue)(seqOp, combOp)	Similar to reduce() but used to return a different type.	rdd.aggregate((0, 0)) ((x, y) => (x._1 + y, x._2 + 1), (x, y) => (x._1 + y._1, x._2 * y._2))	(9, 4)
foreach(func)	Apply the provided function to each element of the RDD.	rdd.foreach(func)	Nothing



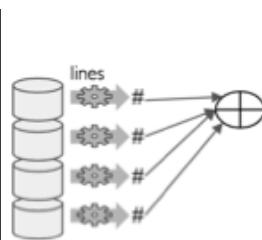
SIT742 | Modern Data Science (G. Li)

99

Caching RDDs

```
lines = sc.textFile("...", 4)
print lines.count()

Count will cause Spark to
• read data
• sum within partitions
• combine sums in driver
```



SIT742 | Modern Data Science (G. Li)

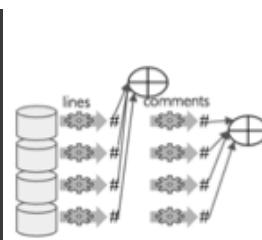
100

Caching RDDs

```
lines = sc.textFile("...", 4)
cmnts = lines.filter(isComment)

print lines.count()
print cmnts.count()

Count will cause Spark to
• read data (again)
• sum within partitions
• combine sums in driver
```



SIT742 | Modern Data Science (G. Li)

101

Caching RDDs

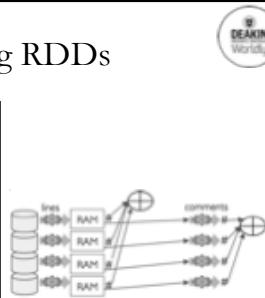
```
lines = sc.textFile("...", 4)

lines.cache()
# save, don't recompute

cmts = lines.filter(isComment)

print lines.count()
print cmts.count()

lines.unpersist()
# remove them from cache
```



102



SIT742 | Modern Data Science (G. Li)

Pair RDDs

- Pair RDD
- Pair RDD Operations
 - Transformations
 - Actions



SIT742 | Modern Data Science (G. Li)

103



Pair RDDs

• Pair RDD

- Similar to MapReduce, Spark supports Key-Value pairs
- Each element of a Pair RDD is a pair tuple
- Allow to act on each key in parallel or regroup data across the network.

```
>>> rdd = sc.parallelize([(1, 2), (3, 4)])
RDD: [(1, 2), (3, 4)]

pairs = lines.map(lambda x: (x.split(" ")[0], x))
# create a pair RDD using the first word as the key in Python
```

104



SIT742 | Modern Data Science (G. Li)

Pair RDD Transformation

• reduceByKey (func)

- Return a new distributed dataset of **(k, v)** pairs where the values for each key are aggregated using the given reduce function **func**, which must be of type **(v, v) → v**

```
>>> rdd = sc.parallelize([(1,2), (3,4), (3,6)])
>>> rdd.reduceByKey(lambda a, b: a + b)
RDD: [(1,2), (3,4), (3,6)] → [(1,2), (3,10)]
```

105



SIT742 | Modern Data Science (G. Li)

Pair RDD Transformation

• sortByKey (func)

- Return a new dataset **(k, v)** pairs sorted by keys in ascending order

```
>>> rdd2 = sc.parallelize([(1,'a'), (2,'c'), (1,'b')])
>>> rdd2.sortByKey()
RDD: [(1,'a'), (2,'c'), (1,'b')] →
[(1,'a'), (1,'b'), (2,'c')]
```

106



SIT742 | Modern Data Science (G. Li)

Pair RDD Transformation

• groupByKey (func)

- Return a new dataset of **(k, iterable<v>)** pairs
- Be careful using this, as it can cause a lot of data movement across the network and create large **iterable** at works

```
>>> rdd2 = sc.parallelize([(1,'a'), (2,'c'), (1,'b')])
>>> rdd2.groupByKey()
RDD: [(1,'a'), (1,'b'), (2,'c')] →
[(1,['a','b']), (2,['c'])]
```

107



SIT742 | Modern Data Science (G. Li)

Pair RDD Transformation

- **mapValues(func)**

- Apply a function to each value of a pair RDD without changing the key

- **keys()**

- Return an RDD of just the keys

```
>>> rdd2 = sc.parallelize([(1,'a'), (2,'c'), (1,'b')])

>>> rdd2.mapValues(lambda x: x+'s')
RDD: [(1,'as'), (2,'cs'), (1,'bs')]

>>> rdd2.keys()
RDD: [1, 2, 1]
```



SIT742 | Modern Data Science (G. Li)

108

Pair RDD Transformation

- **sortByKey()**

- Takes a parameter called **ascending** indicating whether we want it in ascending order
- We can also use our own comparison function

```
rdd.sortByKey(ascending=True, numPartitions=None,
keyfunc = lambda x: str(x))
```



SIT742 | Modern Data Science (G. Li)

109

Pair RDD Transformation

Table 4-1. Transformations on one pair RDD (example: [(1, 2), (3, 4), (3, 6)])

Function name	Purpose	Example	Result
reduceByKey(Func)	Combine values with the same key.	rdd.reduceByKey((x, y) => x + y)	{(1, 2), (2, 3), (3, 10)}
groupByKey()	Group values with the same key.	rdd.groupByKey()	{(1, [2]), (3, [4, 6])}
combineByKey(createCombiner, mergeFunc, mergeCombiner, partitioner)	Combine values with the same key using a different result type.	See Examples 4-12 through 4-14.	
mapValues(Func)	Apply a function to each value of a pair RDD without changing the key.	rdd.mapValues(x => x*1)	{(1, 2), (3, 3), (3, 3)}



SIT742 | Modern Data Science (G. Li)

110

Pair RDD Transformation

Function name	Purpose	Example	Result
values()	Return an RDD of just the values.	rdd.values()	{2, 4, 6}
sortByKey()	Return an RDD sorted by the key.	rdd.sortByKey()	{(1, 2), (3, 4), (3, 6)}
flatMapValues(Func)	Apply a function that returns an iterator to each value of a pair RDD, and for each element returned, produce a key/value entry with the old key. Often used for tokenization.	rdd.flatMapValues(x => x to 5)	{(1, 2), (1, 3), (1, 4), (1, 5), (3, 4), (3, 5)}
keys()	Return an RDD of just the keys.	rdd.keys()	{1, 3, 3}



SIT742 | Modern Data Science (G. Li)

111

Pair RDD Joins

- **X.Join(Y)** is an inner join

- Returns RDD of all pairs of elements with matching keys in X and Y
- Each pair is **(k, (v1, v2))** tuple, where
 - **(k, v1)** is in X
 - **(k, v2)** is in Y

```
>>> x = sc.parallelize([('a', 1), ('b', 4)])
>>> y = sc.parallelize([('a', 2), ('a', 3)])
>>> sorted(x.join(y).collect())
Value: [('a', (1, 2)), ('a', (1, 3))]
```



SIT742 | Modern Data Science (G. Li)

112

Pair RDD Joins

- **X.leftOuterJoin(Y)**

- For each element **(k, v)** in X, resulting RDD will either contain
 - All pairs **(k, (v, w))** for **w** in Y
 - Or pair **(k, (v, None))** if no elements in Y have key **k**

```
>>> x = sc.parallelize([('a', 1), ('b', 4)])
>>> y = sc.parallelize([('a', 2)])
>>> sorted(x.leftOuterJoin(y).collect())
Value: [('a', (1, 2)), ('b', (4, None))]
```



SIT742 | Modern Data Science (G. Li)

113

Pair RDD Joins

- **X.rightOuterJoin(Y)**

- For each element **(k, w)** in Y, resulting RDD will either contain
 - All pairs **(k, (v, w))** for **v** in X
 - Or pair **(k, (None, w))** if no elements in X have key **k**

```
>>> x = sc.parallelize([('a', 1), ('b', 4)])
>>> y = sc.parallelize([('a', 2), ('c', 5)])
>>> sorted(x.rightOuterJoin(y).collect())
Value: [('a', (1, 2)), ('c', (None, 5))]
```



SIT742 | Modern Data Science (G. Li)

114

Pair RDD Joins

- **X.fullOuterJoin(Y)**

- For each element **(k, v)** in X, resulting RDD will either contain all pairs **(k, (v, w))** for **w** in Y, or **(k, (v, None))** if no match in Y have **k**
- For each element **(k, w)** in Y, resulting RDD will either contain all pairs **(k, (v, w))** for **v** in X, or **(k, (None, w))** if no match in X have **k**

```
>>> x = sc.parallelize([('a', 1), ('b', 4)])
>>> y = sc.parallelize([('a', 2), ('c', 5)])
>>> sorted(x.fullOuterJoin(y).collect())
Value: [('a', (1, 2)), ('b', (4, None)), ('c', (None, 5))]
```



SIT742 | Modern Data Science (G. Li)

115

Pair RDD Actions

Table 4-3. Actions on pair RDDs (example: `((1, 2), (3, 4), (3, 6))`)

Function	Description	Example	Result
countByKey()	Count the number of elements for each key.	<code>rdd.countByKey()</code>	<code>{('x', 2), ('y', 2)}</code>
collectAsMap()	Collect the result as a map to provide easy lookup.	<code>rdd.collectAsMap()</code>	<code>Map({('x', 2), ('y', 4), ('z', 6)})</code>
lookup(key)	Return all values associated with the provided key.	<code>rdd.lookup(3)</code>	<code>[4, 6]</code>



SIT742 | Modern Data Science (G. Li)

116

Closures

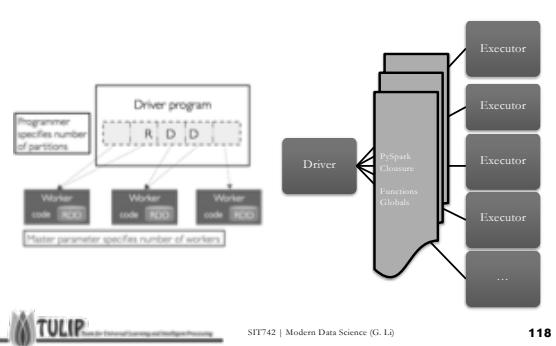
- PySpark Closures
- PySpark Shared Variables
 - Broadcast Variables
 - Accumulators



SIT742 | Modern Data Science (G. Li)

117

PySpark Closures



SIT742 | Modern Data Science (G. Li)

118

PySpark Closures

- Spark automatically creates **closures** for
 - Functions that run on RDDs at workers
 - Any **global variables** used by those workers
- One closure per worker
 - Sent for every task
 - No communication between workers
 - Changes to global variables at workers are not sent to the driver



SIT742 | Modern Data Science (G. Li)

119

PySpark Closures

- Consider the code below (**Local** mode)
 - the code sums the values within the RDD and store it in **counter**, because both the RDD and the variable **counter** are in the same memory space on the driver node.

```
counter = 0
rdd = sc.parallelize(data)

# Wrong: Don't do this!!
rdd.foreach(lambda x: counter += x)

print("Counter value: " + counter)
```



SIT742 | Modern Data Science (G. Li)

120



PySpark Closures

- Consider the code below (**Cluster** Mode)
 - Spark computes/serializes/send the closure to executors
 - Closure contains variables and methods which must be visible for the executor to perform its computations
 - The executors running on separate worker nodes each have their own copy of the closure.

```
counter = 0
rdd = sc.parallelize(data)

# Wrong: Don't do this!!
rdd.foreach(lambda x: counter += x)

print("Counter value: " + counter)
```



SIT742 | Modern Data Science (G. Li)

121



PySpark Closures

- Consider the code below (**Cluster** Mode)
 - The variables within the closure to each executor are copies and thus, when **counter** is referenced within the **foreach** function, it's no longer the **counter** on the driver node. The **counter** in the driver node still there but no longer visible to the executors!

```
counter = 0
rdd = sc.parallelize(data)

# Wrong: Don't do this!!
rdd.foreach(lambda x: counter += x)

print("Counter value: " + counter)
```



SIT742 | Modern Data Science (G. Li)

122



Spark HelloWorld!: Word Count

```
text_file = sc.textFile("hdfs://...")

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

counts.saveAsTextFile("hdfs://...")
```



SIT742 | Modern Data Science (G. Li)

123



PySpark Shared Variables

Broadcast Variables

- keep a read-only variable cached on each machine rather than shipping a copy of it with tasks
- Spark also attempts to distribute broadcast variables using efficient broadcast algorithms to reduce communication cost.



SIT742 | Modern Data Science (G. Li)

124

Accumulators

- Aggregate values from workers back to driver
- Only driver can access value of accumulator
- For tasks, accumulators are write-only



PySpark Shared Variables: Broadcast Variables

- Spark automatically broadcasts the common data needed by tasks within each stage.
 - The data broadcasted this way is cached in serialized form and deserialized before running each task.
- Creating **broadcast variables** is only useful when tasks across multiple stages need the same data or when caching the data in serialized form is important.

```
>>> broadcastVar = sc.broadcast([1, 2, 3])
<pyspark.broadcast.Broadcast object at 0x102789f10>
>>> broadcastVar.value
[1, 2, 3]
```



SIT742 | Modern Data Science (G. Li)

125



PySpark Shared Variables: Broadcast Variables

```
# Lookup the locations of the call signs on the
# RDD contactCounts. We load a list of call sign
# prefixes to country code to support this lookup
signPrefixes = loadCallSignTable()

def processSignCount(sign_count, signPrefixes):
    country = lookupCountry(sign_count[0], signPrefixes)
    count = sign_count[1]
    return (country, count)

countryContactCounts = (contactCounts
    .map(processSignCount)
    .reduceByKey((lambda x, y: x + y)))
```



SIT742 | Modern Data Science (G. Li)

126

PySpark Shared Variables: Broadcast Variables

```
# Lookup the locations of the call signs on the
# RDD contactCounts. We load a list of call sign
# prefixes to country code to support this lookup
signPrefixes = sc.broadcast(loadCallSignTable())

def processSignCount(sign_count, signPrefixes):
    country = lookupCountry(sign_count[0], signPrefixes.value)
    count = sign_count[1]
    return (country, count)

countryContactCounts = (contactCounts
    .map(processSignCount)
    .reduceByKey((lambda x, y: x + y)))
```



SIT742 | Modern Data Science (G. Li)

127

PySpark Shared Variables: Accumulators

- **Accumulator** provides a simple syntax for aggregating values from worker nodes back to the driver program.
- Only driver can read its value
- Used to efficiently implement parallel counters and sums

```
>>> accum = sc.accumulator(0)
>>> rdd = sc.parallelize([1, 2, 3, 4])
>>> def f(x):
...     global accum
...     accum += x
...     rdd.foreach(f)
...     accum.value
Value: 10
```



SIT742 | Modern Data Science (G. Li)

128

PySpark Shared Variables: Accumulators

```
file = sc.textFile(inputFile)
# Create Accumulator[Int] initialized to 0
blankLines = sc.accumulator(0)

def extractCallSigns(line):
    global blankLines # Make the global variable accessible
    if (line == ""):
        blankLines += 1
    return line.split(" ")

callSigns = file.flatMap(extractCallSigns)
print "Blank lines: %d" % blankLines.value
```



SIT742 | Modern Data Science (G. Li)

129

PySpark Shared Variables: Accumulators

- **Accumulators** can be used in actions or transformations:
 - Actions:
 - each task's update to accumulator is applied only once
 - Transformations:
 - no guarantees (use only for debugging)



SIT742 | Modern Data Science (G. Li)

130

Spark File Performance

- Issues to Consider
- File R/W Performance
 - Binary vs. Text
 - Compressed vs. Uncompressed
 - Read vs. Write



SIT742 | Modern Data Science (G. Li)

131

Considerations for File Performance

- Data Sources
 - Where to read data from?
 - Where to save data to?
- Issues to consider
 - Read vs Write Performance
 - Text vs Binary format
 - Uncompressed vs Compressed
 - Environment: Python vs Scala/Java



SIT742 | Modern Data Science (G. Li)

132

File Read/Write Performance

- Test run on 626 MB Text File vs 787 MB Binary File
 - **Python (Pandas)** performance is dependent on library
 - Pandas has no binary file I/O library
 - **Scala/Java:** 6 seconds is the time for sustained read/write. It is often faster due to system catching
 - Read and Write time comparable
 - Binary I/O much faster than Text I/O

	Read Time (Text)	Write Time (Text)	Read Time (Binary)	Write Time (Binary)
Pandas (Python)	36 secs	45 secs	**	**
Scala/Java	18 secs	21 secs	1-6 ⁸ secs	1-6 ⁸ secs



SIT742 | Modern Data Science (G. Li)

133

File Read/Write Performance

- Scala/Java Language
 - Write times much larger than Read
 - Large range of compression time
 - Small range of compressed file sizes



SIT742 | Modern Data Science (G. Li)

134

File Read/Write Performance

- Scala/Java Language
 - Binary I/O still much faster than Text I/O
 - LZ4 Compression ~ Raw I/O Speed

Text File	Read Time	Write Time	File Size
Gzip level 6 (default)	26 secs	98 secs	243 MB
Gzip level 3	25 secs	46 secs	259 MB
Gzip level 1	25 secs	33 secs	281 MB
LZ4 fast	22 secs	24 secs	423 MB
Raw text file	18 secs	21 secs	626 MB



SIT742 | Modern Data Science (G. Li)

135

File Read/Write Performance

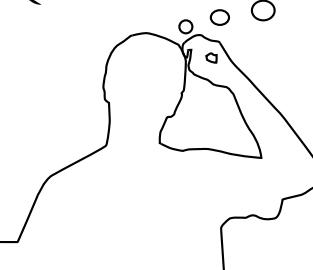
- Summary
 - Uncompressed read and write times are comparable
 - Binary I/O is much faster than Text I/O
 - Compressed reads much faster than compressed writes
 - LZ4 is better than gzip
 - LZ4 compression times approach raw I/O times



SIT742 | Modern Data Science (G. Li)

136

Questions?



SIT742 | Modern Data Science (G. Li)

137

This Session's Readings

- Spark Programming Guide
 - <https://spark.apache.org/docs/latest/programming-guide.html>
- YouTube Video:
 - <https://www.youtube.com/watch?v=xc7Lc8RA8wE&list=tPLxwhieuTaYXmWTBovyyw2NibPfUajk-h4&index=2>

  SIT742 | Modern Data Science (G. Li)  138

Lecture Notes on
Modern Data Science

Module 06: Cloud Computing - SparkSQL

Gang Li
School of Information Technology
Deakin University, VIC 3125, Australia

Unit Learning Outcomes

- ULO1:**
 - Develop knowledge of and discuss new and emerging fields in data science.
 - Cloud Computing and Hadoop/Spark Platform
- ULO4:**
 - Use appropriate platform to collect and process relatively large datasets.
 - Programming in Spark

TULIP | Modern Data Science (G. Li) 1

Road map

```

graph LR
    PS[Programming Spark] --> SSD[SparkSQL and DataFrames]
    SSD --> SS[SparkSQL]
    SSD --> DF[DataFrames]
    SSD --> CDF[Construct DataFrames]
    SSD --> SCV[Spark-CSV]
    SCV --> SC[SQLContext]
    SC --> SS
    SC --> DF
    SC --> CDF
  
```

SIT742 | Modern Data Science (G. Li) 2

Spark SQL and DataFrames

- Tabular Data
- From RDD to DataFrames
- DataFrames

TULIP | Modern Data Science (G. Li) 3

Example of Datasets

visitorid locationid firstaccess lastaccess accesscount macaddress name

visitorid	locationid	firstaccess	lastaccess	accesscount	macaddress	name
1893824	718	23:27.7	56:36.5	38		City Moorabool3 Nth
1893825	718	23:52.1	24:04.0	5		City Moorabool3 Nth
1893826	718	24:05.1	24:05.1	NaN		City Moorabool3 Nth
1893833	412	30:00.2	30:25.2	8	B4:07:F9:B9:F8:61	City Johnstone Pk
1893836	718	39:40.6	53:50.9	88		City Moorabool3 Nth
1893838	718	44:23.0	44:38.1	5.1		City Moorabool3 Nth
1893839	718	44:38.9	44:38.9	1		City Moorabool3 Nth

SIT742 | Modern Data Science (G. Li) 4

Example of Datasets

Family Name	Given Name	Sex	Height	Addr.
Lee	George	M	1.75	VIC
Robins	Tim	M	3.87	SA
Chan	Ruth	F	1.63	NSW
Smith	Jim	M	-	Victoria
Chiang	Yin	F	0.65	TAS
Enders	Robert	M	-	VIC

Attribute: Sex

Values: F, M

What Data Type is this?

How do I make sense of this data?

What are the suitable tools?

What are the possible pitfalls?

TULIP | Modern Data Science (G. Li) 5

Tabular Data

attribute: variables, feature, field



Day, Outlook, Temperature, Humidity, Wind, Play Football

D1	Sunny	Hot	High	Weak	No
D2	Sunny	Hot	High	Strong	No
D3	Overcast	Hot	High	Weak	Yes
D4	Rain	Mild	High	Weak	Yes
D5	Rain	Cool	Normal	Weak	Yes
D6	Rain	Cool	Normal	Strong	No
D7	Overcast	Cool	Normal	Strong	Yes
D8	Sunny	Mild	High	Weak	No
D9	Sunny	Cool	Normal	Weak	Yes
D10	Rain	Mild	Normal	Weak	Yes
D11	Sunny	Mild	Normal	Strong	Yes
D12	Overcast	Mild	High	Weak	Yes
D13	Overcast	Hot	Normal	Weak	Yes
D14	Rain	Mild	High	Strong	No

TULIP: Data for Universal Learning and Intelligent Processing SIT742 | Modern Data Science (G. Li) 6

Tabular Data

- One of the most common data formats
 - CSV, EXCEL, ARFF etc.
- A *table* is a collection of *rows* and *columns*!
 - Each row has an *index* and
 - Each column has a *name*!
- A *cell* is specified by an (*index*, *name*) pair!
 - A cell may or may not have a *value*!
 - A cell's *type* is inferred from its value

TULIP: Data for Universal Learning and Intelligent Processing SIT742 | Modern Data Science (G. Li) 7

Tabular Data

- Limitations
 - Format not well-defined
 - May be missing data values
 - Types may be incorrectly inferred ("2" vs "2.0")
 - No Support for Versioning of Format

TULIP: Data for Universal Learning and Intelligent Processing SIT742 | Modern Data Science (G. Li) 8

Spark Model

- Write programs in terms of transformations on distributed datasets
- Resilient Distributed Datasets** (RDDs)
 - Collections of objects that can be stored in memory or disk across a cluster
 - Parallel functional transformations (map, filter, ...)
 - Automatically rebuilt on failure

TULIP: Data for Universal Learning and Intelligent Processing SIT742 | Modern Data Science (G. Li) 9

Data Challenges and Solutions

Challenges	Solutions
<ul style="list-style-type: none"> Perform ETL to and from various (semi- or unstructured) data sources Perform advanced analytics (e.g. machine learning, graph processing) that are hard to express in relational systems. 	<ul style="list-style-type: none"> SparkSQL and DataFrame API that can perform relational operations on both external data sources and Spark's built-in RDDs. A highly extensible optimizer, Catalyst, that uses features of Scala to add composable rule, control code generation, and define extensions.

TULIP: Data for Universal Learning and Intelligent Processing SIT742 | Modern Data Science (G. Li) 10

Data Characteristic Hierarchy: Data Records



Structured (schema-first)
Semi-Structured (schema-later)
Unstructured (schema-never)

Relational Database
Documents XML
Plain Text
Media
Tagged Text/Media
Formatted Messages

TULIP: Data for Universal Learning and Intelligent Processing SIT742 | Modern Data Science (G. Li) 11

Spark SQL

- **SQL** = *Structured Query Language*
- Some of the Functionality SQL provides:
 - Create, modify, delete relations
 - Add, modify, remove tuples
 - Specify queries to find tuples matching criteria

```
Select * from Users Where Users.age < 21
Select S.name, E.UID
From Students S, EnrolledUnits E
Where S.SID = E.SID
```

TULIP Data for Unstructured Learning and Intelligent Processing SIT742 | Modern Data Science (G. Li) 12

Spark SQL

- **Spark SQL** allows you to manipulate distributed data with SQL queries. Currently, two SQL dialects are supported.
 - If you're using a Spark **SQLContext**, the only supported dialect is "**sql**", a rich subset of SQL 92.
 - If you're using a **HiveContext**, the default dialect is "**hiveql**", corresponding to Hive's SQL dialect.
 - "**sql**" is also available, but "**hiveql**" is a richer dialect.



TULIP Data for Unstructured Learning and Intelligent Processing SIT742 | Modern Data Science (G. Li) 13

Spark SQL

- You issue SQL queries through a **SQLContext** or **HiveContext**, using the **sql()** method which returns a **DataFrame**.
 - You can mix DataFrame methods and SQL queries in the same code.
- To use SQL, you *must* either:
 - query a persisted Hive table, or
 - make a *table alias* for a DataFrame, using **registerTempTable()**

TULIP Data for Unstructured Learning and Intelligent Processing SIT742 | Modern Data Science (G. Li) 14

DataFrames

- **DataFrame** is a size-mutable, potentially heterogeneous tabular data structure with labeled axes (i.e., rows and columns).
 - a recent extension to **Spark RDD API** (early 2015)
 - *Distributed* collection of data organized into named columns
 - Equivalent to **Python/Pandas** and **R DataFrame**, but distributed and with richer optimizations under the hood
 - designed from the ground-up to support modern big data and data science applications
 - Types of columns inferred from values

TULIP Data for Unstructured Learning and Intelligent Processing SIT742 | Modern Data Science (G. Li) 15

DataFrames

- **DataFrame** has the following features:
 - Ability to scale from kilobytes of data on a single laptop to petabytes on a large cluster
 - Support for a wide array of data formats and storage systems
 - State-of-the-art optimization and code generation through the Spark SQL **Catalyst** optimizer
 - Seamless integration with all big data tooling and infrastructure via Spark
 - APIs for Python, Java, Scala, and R (in development via **SparkR**)

TULIP Data for Unstructured Learning and Intelligent Processing SIT742 | Modern Data Science (G. Li) 16

DataFrames

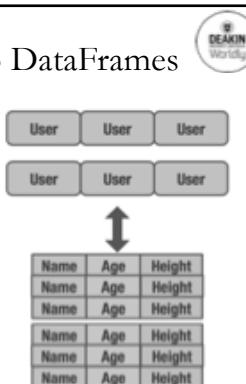
```

graph TD
    JDBC[ JDBC ] --> SparkSQL[ Spark SQL ]
    Console[ Console ] --> SparkSQL
    UserPrograms[ User Programs (Java, Scala, Python) ] --> SparkSQL
    CatalystOptimizer[ Catalyst Optimizer ] --> SparkSQL
    CatalystOptimizer --> RDDS[ Resilient Distributed Datasets ]
  
```

TULIP Data for Unstructured Learning and Intelligent Processing SIT742 | Modern Data Science (G. Li) 17

From RDDs to DataFrames

- RDDs**
 - Functional transformations on partitioned collections of opaque objects
- DataFrames**
 - Declarative transformations on partitioned collections of tuples



SIT742 | Modern Data Science (G. Li)

18

From RDDs to DataFrames

- DataFrames** are built on top of the Spark RDD API.
 - You can use normal RDD operations on DataFrames.
 - Stick with the DataFrame API, wherever possible.
 - Using RDD operations will often give you back an RDD, not a DataFrame.
 - The DataFrame API is likely to be more efficient, because it can optimize the underlying operations with Catalyst.

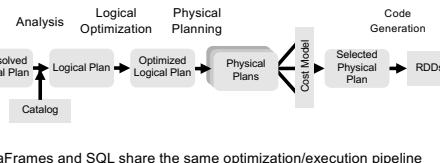
TULIP: Towards Universal Learning and Intelligent Processing

SIT742 | Modern Data Science (G. Li)

19

From RDDs to DataFrames

- Plan Optimization & Execution
 - Represented internally as a “*logical plan*”
 - Execution is lazy, allowing it to be optimized by Catalyst



DataFrames and SQL share the same optimization/execution pipeline

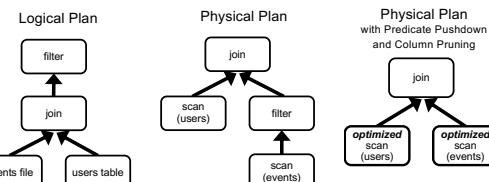
SIT742 | Modern Data Science (G. Li)

20

From RDDs to DataFrames

- Plan Optimization & Execution

```
joined = users.join(events, users.id == events.uid)
filtered = joined.filter(events.date >= "2015-01-01")
```



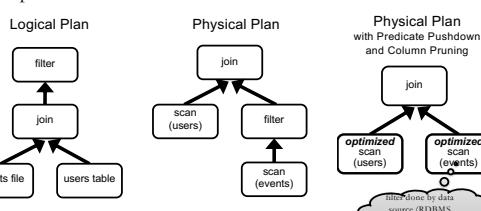
TULIP: Towards Universal Learning and Intelligent Processing

SIT742 | Modern Data Science (G. Li)

21

From RDDs to DataFrames

- Plan Optimization & Execution
 - Data Sources API can automatically prune columns and push filters to the source

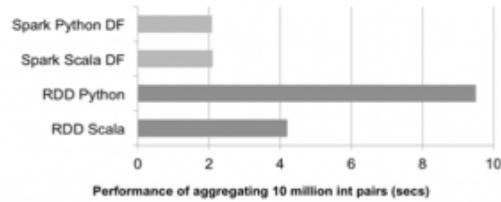


SIT742 | Modern Data Science (G. Li)

22

From RDDs to DataFrames

- DataFrames can be significantly faster than RDDs
- Performance is irrelevant to language



TULIP: Towards Universal Learning and Intelligent Processing

23

Data Sources Supported by DataFrames

- DataFrame execution is done internally with RDDs making interoperation with outside sources easy.
- Easy loading/saving DataFrames

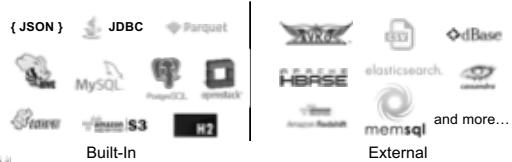


SIT742 | Modern Data Science (G. Li)

24

Data Sources Supported by DataFrames

- Efficient data access by Spark SQL query optimizer
 - Have interfaces allowing optimizations to be pushed down to data source
 - e.g. avoid reading unnecessary data for a query



SIT742 | Modern Data Science (G. Li)

25

SparkSQL and DataFrames

- SparkSession
- Construct DataFrame
 - Table
 - File
 - RDD
- Spark-CSV



SIT742 | Modern Data Science (G. Li)

26

SparkSession

- Entry to all relational functionality in Spark is **SparkSession** class or its decedents.

```
from pyspark.sql import SparkSession
sqlContext = SparkSession \
.builder \
.appName("Python Spark SQL basic example") \
.config("spark.some.config.option", "some-value") \
.getOrCreate()
```

SIT742 | Modern Data Science (G. Li)

27

Construct a DataFrame from Tables

- **DataFrames** can be constructed from a wide array of sources such as:
 - tables in Hive, external databases, structured data files, or existing RDDs.
 - The example shows creation from a Hive Table

```
#Construct a DataFrame from a "users" table in Hive.
df = sqlContext.table("users")

# Displays the content of the DataFrame to stdout
df.show()
```

SIT742 | Modern Data Science (G. Li)

28

Construct a DataFrame from Data with Schema

- The (changing) role of Schema
 - **Schema** specify the **structure** and **types** of a data repository, e.g. the types of each column in a table.
 - They may also specify constraints **within** or **between** data fields.
 - Traditional databases are **schema-on-write**.
 - You cannot load data into a table without a schema.
 - Newer (NoSQL) data stores are **schema-on-read** or **schemaless**:
 - You can defer applying a schema until you read the data, or avoid schema altogether.



SIT742 | Modern Data Science (G. Li)

29

Construct a DataFrame from Data with Schema

- We can easily create **DataFrames** from those data formats that have a schema, which includes the column names and data types of **DataFrames**.
 - A **Parquet** table has a schema that Spark can use, it also allows Spark to be efficient about how it parses data
 - The example shows creation from a Parquet file

```
# Construct a DataFrame from a PARQUET file in installed folder.
df = sqlContext.read.load("examples/src/main/resources/users.parquet")

# Displays the content of the DataFrame to stdout
df.show()
```

SIT742 | Modern Data Science (G. Li)

30

Construct a DataFrame from Data with Schema

- We can easily create **DataFrames** from those data formats that have a schema, which includes the column names and data types of **DataFrames**.
 - The example shows creation from a JSON file

```
# Construct a DataFrame from a JASON file in installed folder.
df = sqlContext.read.json("examples/src/main/resources/people.json")

# Displays the content of the DataFrame to stdout
df.show()
```

SIT742 | Modern Data Science (G. Li)

31

Construct a DataFrame from RDD

- Spark SQL can convert an **RDD** of **Row** objects to a **DataFrame**, inferring the data types.
 - Rows** are constructed by passing a list of key/value pairs as kwargs to the **Row** class.
 - The keys of this list define the column names of the table, and
 - The types are inferred by looking at the first row.

TULIP: Tools for Universal Learning and Intelligent Processing

SIT742 | Modern Data Science (G. Li)

32

Construct a DataFrame from RDD (1)

```
# sc is an existing SparkContext.
from pyspark.sql import SQLContext, Row
sqlContext = SQLContext(sc)

# Load a text file and convert each line to a Row.
lines = sc.textFile("examples/src/main/resources/people.txt")
parts = lines.map(lambda l: l.split(","))
people = parts.map(lambda p: Row(name=p[0], age=int(p[1])))

# Infer the schema, and register the DataFrame as a table.
schemaPeople = sqlContext.createDataFrame(people)

schemaPeople.registerTempTable("people")
```

SIT742 | Modern Data Science (G. Li)

33

Construct a DataFrame from RDD (2)

```
# Import SQLContext and data types
from pyspark.sql import SQLContext
from pyspark.sql.types import *

# sc is an existing SparkContext.
sqlContext = SQLContext(sc)
# Load a text file and convert each line to a tuple.
lines = sc.textFile("examples/src/main/resources/people.txt")
parts = lines.map(lambda l: l.split(","))
people = parts.map(lambda p: (p[0], p[1].strip()))

# The schema is encoded in a string.
schemaString = "name age"
fields = [StructField(field_name, StringType(), True) for field_name in schemaString.split()]
schema = StructType(fields)
# Apply the schema to the RDD.
schemaPeople = sqlContext.createDataFrame(people, schema)
# Register the DataFrame as a table.
schemaPeople.registerTempTable("people")
```

SIT742 | Modern Data Science (G. Li)

34

Example: Construct a DataFrame from CSV

- Suppose we have a text file that has no schema
 - We can infer the schema
- | | |
|---|---|
| First Name: string
Last Name: string
Gender: string
Age: integer | Rosalita,Ramirez,F,14
Ally,Garcia,F,39
Claire,McBride,F,23
Abigail,Cottrell,F,75
José,Rivera,M,59
Ravi,Dasgupta,M,25 |
|---|---|

SIT742 | Modern Data Science (G. Li)

35

Example: Construct a DataFrame from CSV (1)

```
from pyspark.sql import Row

rdd = sc.textFile("examples/src/main/resources/people.csv")
Person = Row('first_name', 'last_name', 'gender', 'age')

def line_to_person(line):
    cells = line.split(",")
    cells[3] = int(cells[3])
    return Person(*cells)

peopleRDD = rdd.map(line_to_person)

df = peopleRDD.toDF()
# DataFrame[first_name: string, last_name: string, gender: string, age: bigint]
```



SIT742 | Modern Data Science (G. Li)

36

Example: Construct a DataFrame from CSV (2)

```
rdd = sc.textFile("examples/src/main/resources/people.csv")

def line_to_person(line):
    cells = line.split(",")
    return tuple(cells[0:3] + [int(cells[3])])

peopleRDD = rdd.map(line_to_person)

df = peopleRDD.toDF(("first_name", "last_name", "gender", "age"))
```



SIT742 | Modern Data Science (G. Li)

37

Extensions of DataFrames API

- **DataFrames** API can be extended to understand additional input formats (or, input *sources*).
 - For instance, if you're dealing with CSV files, a *very* common data file format, you can use the *spark-csv* package
 - This package augments the DataFrames API so that it understands CSV files.



SIT742 | Modern Data Science (G. Li)

38

Spark-CSV

- Let us assume our data file has a header

```
firstname,lastname,gender,age
Rosalita,Ramirez,F,14
Ally,Garcia,F,39
Claire,McBride,F,23
Abigail,Cottrell,F,75
José,Rivera,M,59
Ravi,Dasgupta,M,25
```



SIT742 | Modern Data Science (G. Li)

39

Spark-CSV

- With *spark-csv*, we can simply create a **DataFrame** directly from our CSV file.
 - *spark-csv* uses the header to infer the schema, but the column types will always be *string*.

```
# Python
df = sqlContext.read.format("com.databricks.spark.csv").\
    load("people.csv", header="true")

// df: org.apache.spark.sql.DataFrame = [fname: string,
// lname: string, gender: string, age: string]
```



SIT742 | Modern Data Science (G. Li)

40

Spark-CSV

- You can also declare the schema programmatically, which allows you to specify the column types.

```
from pyspark.sql.types import *
schema = StructType([StructField("firstName", StringType(), False),
                     StructField("lastName", StringType(), False),
                     StructField("gender", StringType(), False),
                     StructField("age", IntegerType(), False)])
df = sqlContext.read.format("com.databricks.spark.csv").\
    schema(schema).\
    load("people.csv")
```



SIT742 | Modern Data Science (G. Li)

41

DataFrame Operations



- Transformations
- Actions

TULIP Data for Universal Learning and Intelligent Processing

SIT742 | Modern Data Science (G. Li)

42

DataFrame Operations

- DataFrames are *lazy*
 - Transformations** contribute to the query plan, but they don't execute anything.
 - filter, select, drop, intersect, join etc.
 - Actions** cause the execution of the query
 - Spark initiates a distributed read of the data source
 - The data flows through the transformations (the RDDs resulting from the Catalyst query plan)
 - The result of the action is pulled back into the driver JVM.

TULIP Data for Universal Learning and Intelligent Processing

SIT742 | Modern Data Science (G. Li)

43

DataFrames Operations



A distributed collection of data presented in named columns. A DataFrame is equivalent to a relational table in Spark SQL, and can be created using various functions or applications.

```
people = spark.read.json("people.json")
```

Once created, it can be manipulated using the various domain-specific language (DSL) functions defined in `DataFrame`. To select a column from the data frame, use the `select` method.

```
people.select("name").show()
```

A more concrete example:

```
# To create DataFrame using RDDs
people = sqlContext.createDataFrame(rdd, ["name", "age", "dept"])
people.rdd.getNumPartitions() # 2
people.rdd.map(lambda x: (x[0], x[1] + 10)).map(lambda x: (x[0], x[1] - department_id))
people.show()
```

New in version 1.3.

```
avg("age")
# Aggregate on the entire DataFrame without group (similar to SQL's group_by())
avg("age").show()
```

SIT742 | Modern Data Science (G. Li)

44

DataFrames Operations

- Columns**
 - A DataFrame *column* is an abstraction. It provides a common column-oriented view of the underlying data, *regardless* of how the data is really organized.
 - A place (*a cell*) for a data value to reside, within a row of data. This cell can have several states:
 - empty (null), missing (not there at all), or contains a (typed) value (non-null)
 - A collection of those cells, from multiple rows

TULIP Data for Universal Learning and Intelligent Processing

SIT742 | Modern Data Science (G. Li)

45

DataFrames Operations

- Columns**
 - Assume we have a **DataFrame**, `df`, that reads a data source with “`first`”, “`last`”, and “`age`” columns.
 - In Python, it’s possible to access a **DataFrame**’s columns either by attribute (`df.age`) or by indexing (`df['age']`).

```
# Python
df = sqlContext.read.format("com.databricks.spark.csv").\
    load("people.csv", headers="true")

// df: org.apache.spark.sql.DataFrame = [fname: string,
// lname: string, gender: string, age: string]
```

TULIP Data for Universal Learning and Intelligent Processing

SIT742 | Modern Data Science (G. Li)

46

DataFrames Operations

- Show()**
 - You can look at the first *n* elements in a **DataFrame** with the `show()` method.
 - If not specified, *n* defaults to 20.
 - This method is an *action*: It:
 - reads (or re-reads) the input source
 - executes the RDD DAG across the cluster
 - pulls the *n* elements back to the driver JVM
 - displays those elements in a tabular form

TULIP Data for Universal Learning and Intelligent Processing

SIT742 | Modern Data Science (G. Li)

47

DataFrames Operations

- Select()**

- behaves like a SQL SELECT, allowing you to limit the results to specific columns.
- also allows you create on-the-fly *derived columns*.

```
df.select(df['first_name'], df['age'], df['age'] > 49).show(5)

+-----+-----+
|first_name|age|
+-----+-----+
|  Erin| 42|   false|
| Claire| 23|   false|
| Norman| 81|    true|
| Miguel| 64|    true|
+-----+-----+
```



SIT742 | Modern Data Science (G. Li)

48

DataFrames Operations

- Select()**

- you can also use SQL

```
df.registerTempTable("names")
sqlContext.sql("SELECT first_name, age, age > 49 FROM names").\
show(5)

+-----+-----+
|first_name|age|
+-----+-----+
|  Norman| 81|    true|
| Miguel| 64|    true|
+-----+-----+
```



SIT742 | Modern Data Science (G. Li)

49

DataFrames Operations

- filter()**

- allows you to filter rows out of your results.

```
df.filter(df['age'] > 49).\
select(df['first_name'], df['age']).\
show()

+-----+-----+
|firstName|age|
+-----+-----+
|  Norman| 81|
| Miguel| 64|
| Abigail| 75|
+-----+-----+
```



SIT742 | Modern Data Science (G. Li)

50

DataFrames Operations

- filter()**

- allows you to filter rows out of your results.
- Here's the SQL version.

```
sqlContext.sql("SELECT first_name, age FROM names " + \
               "WHERE age > 49").show()

+-----+-----+
|firstName|age|
+-----+-----+
|  Norman| 81|
| Miguel| 64|
| Abigail| 75|
+-----+-----+
```



SIT742 | Modern Data Science (G. Li)

51

DataFrames Operations

- orderby()**

- allows you to sort your results.

```
df.filter(df['age'] > 49).\
select(df['first_name'], df['age']).\
orderBy(df['age'].desc(), df['first_name']).show()

+-----+-----+
|firstName|age|
+-----+-----+
|  Norman| 81|
| Abigail| 75|
| Miguel| 64|
+-----+-----+
```



SIT742 | Modern Data Science (G. Li)

52

DataFrames Operations

- orderby()**

- allows you to sort your results.
- Here's the SQL version.

```
sqlContext.sql("SELECT first_name, age FROM names " + \
               "WHERE age > 49 ORDER BY age DESC, first_name").show()

+-----+-----+
|firstName|age|
+-----+-----+
|  Norman| 81|
| Miguel| 64|
| Abigail| 75|
+-----+-----+
```



SIT742 | Modern Data Science (G. Li)

53

DataFrames Operations

- **groupBy()**

- Often used with **count()**, **groupBy()** groups data items by a specific column value

```
df.groupBy("age").count().show()
```

age	count
39	1
42	2
64	1
75	1
81	1

SIT742 | Modern Data Science (G. Li)

54

DataFrames Operations

- **groupBy()**

- and SQL is not surprising

```
SQLContext.sql("SELECT age, count(age) FROM names " + \
               "GROUP BY age").show()
```

age	count
39	1
42	2
64	1
75	1
81	1

SIT742 | Modern Data Science (G. Li)

55

DataFrames Operations

- **alias()**

- allows you to rename a column.
- It's especially useful with generated columns.

```
df.select(df['first_name'], df['age'], \
          (df['age'] < 30).alias('young')).show(5)
```

first_name	age	young
Erin	42	false
Claire	23	true
Norman	81	false
Miguel	64	false
Rosalita	14	true

SIT742 | Modern Data Science (G. Li)

56

DataFrames Operations

- **alias()**

- and of course, in SQL

```
SQLContext.sql("SELECT firstName, age, age<30 AS young " + \
               "FROM names").show()
```

first_name	age	young
Erin	42	false
Claire	23	true
Norman	81	false
Miguel	64	false
Rosalita	14	true

SIT742 | Modern Data Science (G. Li)

57

DataFrames Operations

- **join()**

- Let's assume we have a second file, a JSON file that contains records like this:

```
{
  "firstName": "Erin",
  "lastName": "Shannon",
  "medium": "oil on canvas"
},
{
  "firstName": "Norman",
  "lastName": "Lockwood",
  "medium": "metal
(sculpture)"
},
```

SIT742 | Modern Data Science (G. Li)

58

DataFrames Operations

- **join()**

- We can load and then join a second DataFrame

```
df2='sqlContext.read.json("artists.json")' #Schema inferred as DataFrame(firstName: string, lastName: string, medium: string)
df.join(
  df2,
  df.firstName == df2.firstName and df.lastName == df2.lastName
).show()
```

firstName	lastName	gender	age	medium
Norman	Lockwood	M	81	oil on canvas
Erin	Shannon	F	42	oil on canvas
Rosalita	Ramirez	F	14	charcoal
Miguel	Ruiz	M	64	oil on canvas

SIT742 | Modern Data Science (G. Li)

59

DataFrames Operations

- join()**

- We can also select some of the columns

```
df2 = sqlContext.read.json("artists.json")
# schemaInferred="DataFrame[firstName: string,lastName: string,medium: string]
df = df2
df.firstName == df2.firstName and df.lastName == df2.lastName )
df3.select("firstName", "lastName", "age", "medium").show()

+-----+-----+-----+
|firstName|lastName|age|medium|
+-----+-----+-----+
| Norman| Lockwood| 81| metal (sculpture)|
|   Erin| Shannon| 42| oil on canvas|
| Rosalita| Ramirez| 14| charcoal|
| Miguel| Ruiz| 64| oil on canvas|
+-----+-----+-----+
```



SIT742 | Modern Data Science (G. Li)

60

Example of DataFrames Operations

```
from pyspark.sql import SQLContext
from pyspark.sql.types import *
# sc is an existing SparkContext.
sqlContext = SQLContext(sc)
# Load a text file and convert each line to a tuple.
lines = sc.textFile("examples/src/main/resources/people.txt")
parts = lines.map(lambda l: l.split(","))
people = parts.map(lambda p: (p[0], p[1].strip()))
# The schema is encoded in a string.
schemaString = "name age"
fields = [StructField(field_name, StringType(), True) for field_name in \
          schemaString.split()]
schema = StructType(fields)
# Apply the schema to the RDD.
schemaPeople = sqlContext.createDataFrame(people, schema)
# Register the DataFrame as a table.
schemaPeople.registerTempTable("people")
# SQL can be run over DataFrames that have been registered as a table.
results = sqlContext.sql("SELECT name FROM people")
# The results of SQL queries are RDDs and support all the normal RDD operations.
names = results.map(lambda p: "Name: " + p.name)
for name in names.collect(): print(name)
```



SIT742 | Modern Data Science (G. Li)

61

Writing DataFrames



- to Files
- to HIVE Tables
- Save Modes



SIT742 | Modern Data Science (G. Li)

62

Writing DataFrames into Files

- You can write DataFrames out, as well.
 - When doing ETL, this is a *very* common requirement.
 - In most cases, if you can read a data format, you can write that data format, as well.
 - If you're writing to a text file format (e.g., JSON), you'll typically get multiple output files.

```
df.write.format("json").save("/path/to/directory")
df.write.format("parquet").save("/path/to/directory")
```



SIT742 | Modern Data Science (G. Li)

63

Writing DataFrames into Hive

- When with a HiveContext, you can save a DataFrame as a persistent table, with the `saveAsTable()` method.
 - Unlike `registerTempTable()`, `saveAsTable()` materializes the DataFrame (i.e., runs the DAG) and creates a pointer to the data in the Hive metastore.
 - Persistent tables will exist even after your Spark program has restarted.
 - By default, `saveAsTable()` will create a *managed table*.
 - the metastore controls the location of the data.
 - Data in a managed table is also deleted automatically when the table is dropped.



SIT742 | Modern Data Science (G. Li)

64

Writing DataFrames: Save Modes

Scala/Java	Python	Meaning
<code>SaveMode.ErrorIfExists</code> (default)	"error"	if output data or table already exists, an exception is expected to be thrown.
<code>SaveMode.Append</code>	"append"	if output data or table already exists, append contents of the DataFrame to existing data.
<code>SaveMode.Overwrite</code>	"overwrite"	if output data or table already exists, replace existing data with contents of DataFrame.
<code>SaveMode.Ignore</code>	"ignore"	if output data or table already exists, do not write DataFrame at all.



SIT742 | Modern Data Science (G. Li)

65

Writing DataFrames: Save Modes

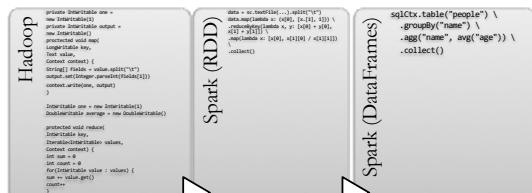
- These save modes do *not* utilize any locking and are *not* atomic.
 - Thus, it is *not* safe to have multiple writers attempting to write to the same location.
 - Additionally, when performing an overwrite, the data will be deleted before writing out the new data.



SIT742 | Modern Data Science (G. Li)

66

MapReduce → RDD → DataFrames → ...



SIT742 | Modern Data Science (G. Li)

67

What is Next? Spark Datasets!

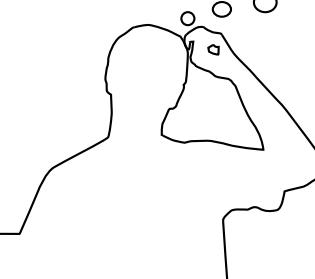
- Spark Datasets, an extension of the DataFrame API that provides a *type-safe, object-oriented programming interface*.
 - Spark 1.6 and later versions includes an API preview of Datasets, but no Python API in the current release
 - Like DataFrames, Datasets take advantage of Spark's Catalyst optimizer by exposing expressions and data fields to a query planner.
 - A Dataset is a strongly-typed, immutable collection of objects that are mapped to a relational schema.



SIT742 | Modern Data Science (G. Li)

68

Questions?



SIT742 | Modern Data Science (G. Li)

69

This Session's Readings

- Spark SQL, DataFrames and Datasets
 - <http://spark.apache.org/docs/latest/sql-programming-guide.html>
 - <https://databricks.com/blog/2015/02/17/introducing-dataframes-in-spark-for-large-scale-data-science.html>
 - <http://spark.apache.org/docs/latest/sql-programming-guide.html#starting-point-sparksession>
- Usage
 - <https://databricks.com/blog/2015/06/02/statistical-and-mathematical-functions-with-dataframes-in-spark.html>



SIT742 | Modern Data Science (G. Li)

70

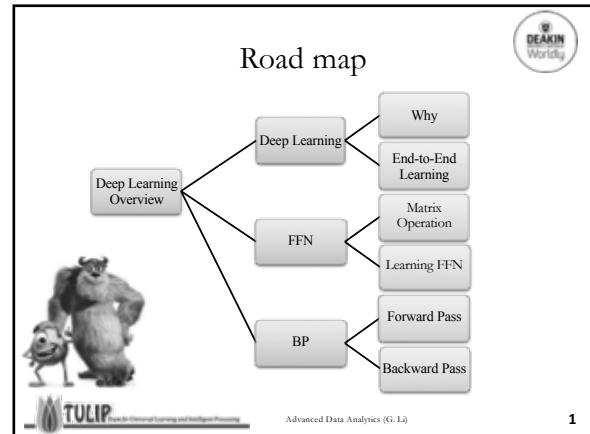
Lecture Notes on
Advanced Data Analytics

Module 05: Deep Learning Overview

Gang Li
School of Information Technology
Deakin University, VIC 3125, Australia

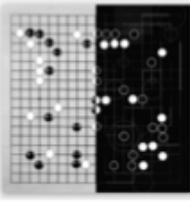


TULIP Team for Unsupervised Learning and Intelligent Processing



Deep Learning

AlphaGo



TULIP Team for Unsupervised Learning and Intelligent Processing

Advanced Data Analytics (G. Li)

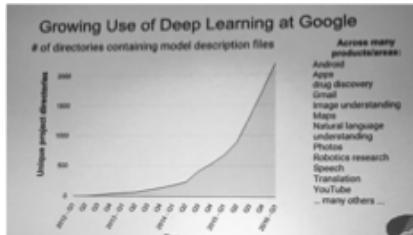
2



Deep learning trends at Google

- SIGMOD 2016, Jeff Dean

Growing Use of Deep Learning at Google



Advanced Data Analytics (G. Li)

4

- ## What is Deep Learning?
- Fast answer:
 - Fast answer: multilayer perceptron (aka deep neural networks) of the 1980s rebranded in 2006.
 - But early nets go stuck at 1-2 hidden layers.
 - Slow answer:
 - distributed representation, multiple steps of computation, modelling the compositionality of the world, a better prior, advances in compute, data & optimization, neural architectures, etc.
- 
- Advanced Data Analytics (G. Li)
- 5

From ANN to Deep Learning

- 1958: Perceptron (linear model)
- 1969: Perceptron has limitation
- 1980s: Multi-layer perceptron
 - Do not have significant difference from DNN today
- 1986: Backpropagation
 - Usually more than 3 hidden layers is not helpful
- 1989: 1 hidden layer is “good enough”
 - why deep?



Advanced Data Analytics (G. Li)

6

From ANN to Deep Learning

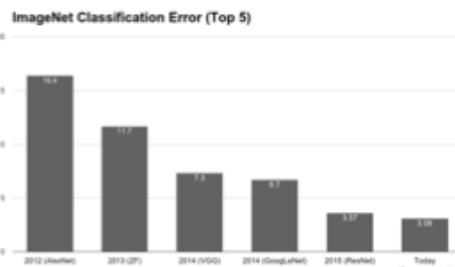
- 2006: RBM initialization
- 2009: GPU
- 2012: win ILSVRC image competition
- 2015.2:
 - Image recognition surpassing human-level performance
- 2016.3: Alpha GO beats Lee Sedol
- 2016.10:
 - Many system as good as humans



Advanced Data Analytics (G. Li)

7

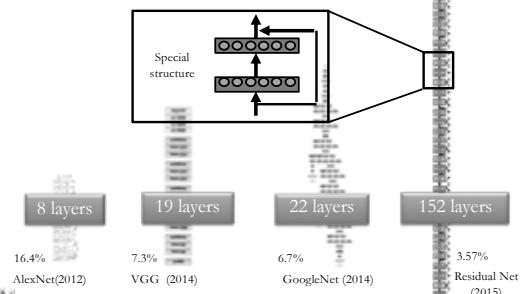
Deep Works



Advanced Data Analytics (G. Li)

8

Deep Works



Advanced Data Analytics (G. Li)

9

The Best of Machine Learning

- Strong/flexible priors:
 - Good features: Feature engineering
 - Data structure: HMM, CRF, MRF, Bayesian nets
 - Model structure, VC-dimension, regularization, sparsity:
 - SVM, compressed sensing
 - Manifold assumption, class/region separation:
 - Metric + semi-supervised learning
 - Factors of variation: PCA, ICA, FA
- Uncertainty quantification:
 - Bayesian, ensemble: RF, GBM
- Sharing statistical strength:
 - model reuse: transfer learning, domain adaption, multitask learning, lifelong learning

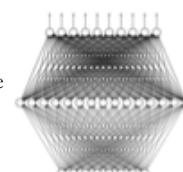


Advanced Data Analytics (G. Li)

10

Universality Theorem

- Shallow network can represent any function.
- Given **enough** hidden neurons, Any continuous function $f: R^N \rightarrow R^M$ can be realized by a network with one hidden layer
- However, using deep structure is more effective.



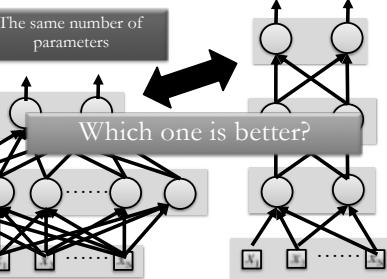
Reference for the reason:
<http://neuralnetworksanddeeplearning.com/chap4.html>



Advanced Data Analytics (G. Li)

11

Fat + Short v.s. Thin + Tall



12



Advanced Data Analytics (G. Li)

Fat + Short v.s. Thin + Tall

Layer X Size	Word Error Rate (%)	Layer X Size	Word Error Rate (%)
1 X 2k	24.2		
2 X 2k	20.4		
3 X 2k	18.4		
4 X 2k	17.8		
5 X 2k	17.2	1 X 3772	22.5
7 X 2k	17.1	1 X 4634	22.6
		1 X 16k	22.1

• Seide, Frank, Gang Li, and Dong Yu. "Conversational Speech Transcription Using Context-Dependent Deep Neural Networks." *Interspeech*. 2011.

Advanced Data Analytics (G. Li)

13

Modularization

- Deep → Modularization
 - Don't put everything in your main function.



<http://riuboney.github.io/2015/10/18/theoretical-motivations-deep-learning.html>



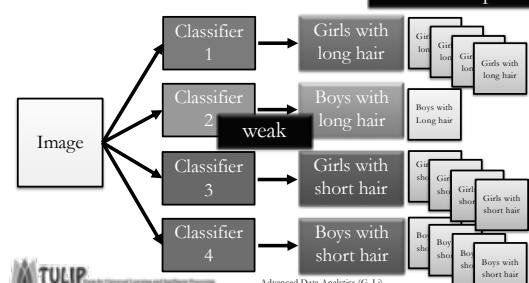
Advanced Data Analytics (G. Li)

14

Modularization

- Deep → Modularization

Little examples

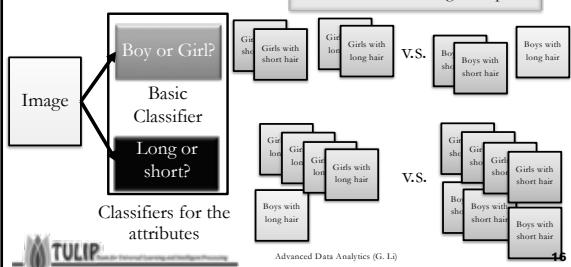


Advanced Data Analytics (G. Li)

15

Modularization

- Deep → Modularization
 - Each basic classifier can have sufficient training examples.



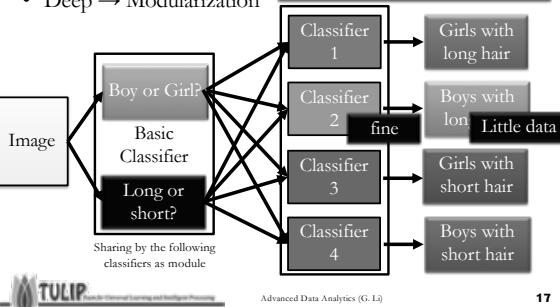
Advanced Data Analytics (G. Li)

16

Modularization

- Deep → Modularization

can be trained by little data



Advanced Data Analytics (G. Li)

17

Modularization

- Deep → Modularization → Less training data?

The most basic classifiers Use 1st layer as module to build classifiers Use 2nd layer as module

The modularization is automatically learned from data.

18

TULIP: Tools for Unsupervised Learning and Intelligent Processing

Advanced Data Analytics (G. Li)

Modularization

- Deep → Modularization

The most basic classifiers Use 1st layer as module to build classifiers Use 2nd layer as module

19

TULIP: Tools for Unsupervised Learning and Intelligent Processing

Reference: Zeiler, M. D., & Fergus, R. (2014). Visualizing and understanding convolutional networks. In *Computer Vision—ECCV 2014* (pp. 818–833).

Advanced Data Analytics (G. Li)

End-to-end Learning

- End-to-end training:
 - What each function should do is learned automatically

Model Hypothesis Functions

Simple Function 1 Simple Function 2 ... Simple Function N

A very complex function

"Hello"

20

TULIP: Tools for Unsupervised Learning and Intelligent Processing

Advanced Data Analytics (G. Li)

End-to-end Learning

- Shallow Approach
 - Each box is a simple function in the production line:

monkey?

feature extr.

pooling

encoding

classification

hand-crafted

learned from data

21

TULIP: Tools for Unsupervised Learning and Intelligent Processing

Advanced Data Analytics (G. Li)

End-to-end Learning

- Deep Learning
 - All functions are learned from data

f₁ **f₂** **f₃** **f₄**

"monkey"

22

TULIP: Tools for Unsupervised Learning and Intelligent Processing

Advanced Data Analytics (G. Li)

To learn more ...

- Do Deep Nets Really Need To Be Deep? (by Rich Caruana)
 - <http://research.microsoft.com/apps/video/default.aspx?id=232373&cr=1>

Do deep nets really need to be deep?

Rich Caruana
Microsoft Research
University of Toronto
UofT Interdisciplinary Data Science Program

Thanks also to: Oleg Vinyals, Armand Joulin, Sumit Majumdar, Mohammad Aljemai, Ali Farhadi, Jia Deng, and Wei Dong

Yes!

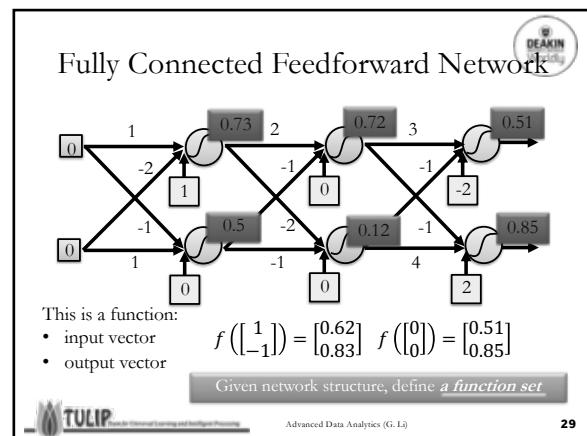
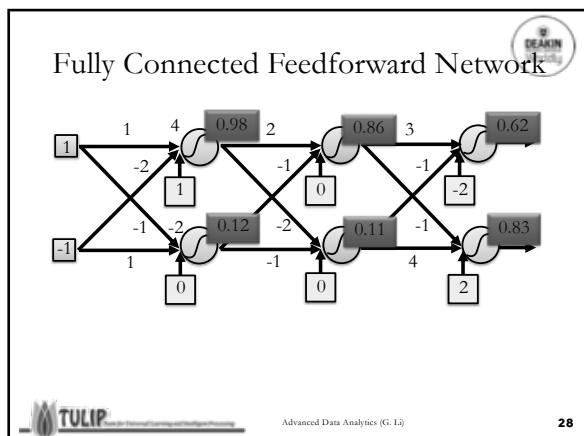
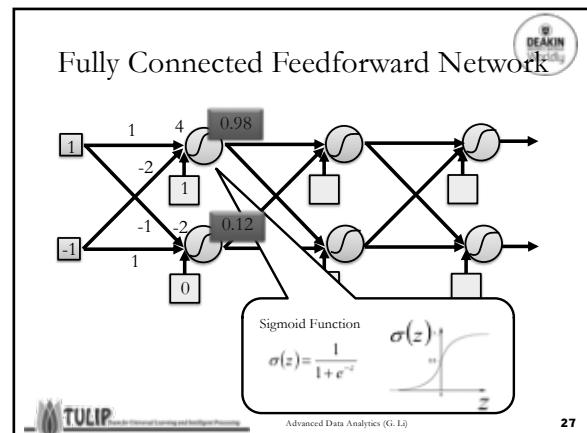
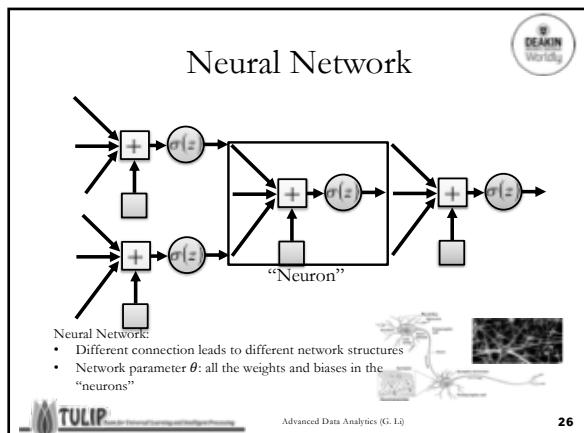
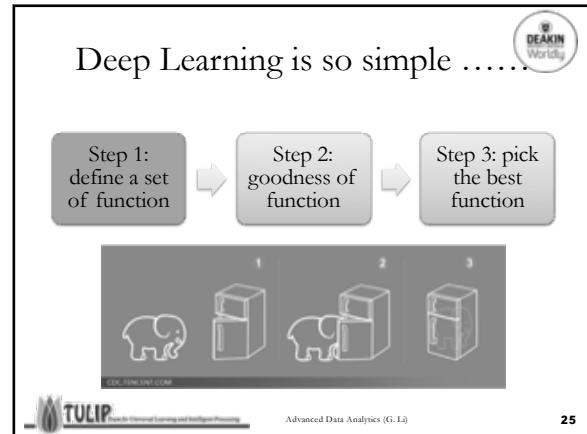
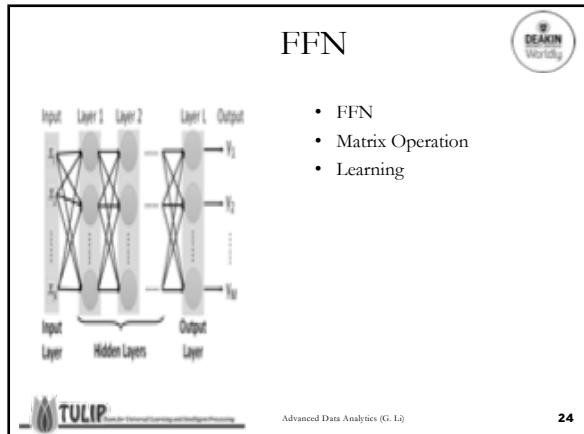
Thank You
Any Questions?

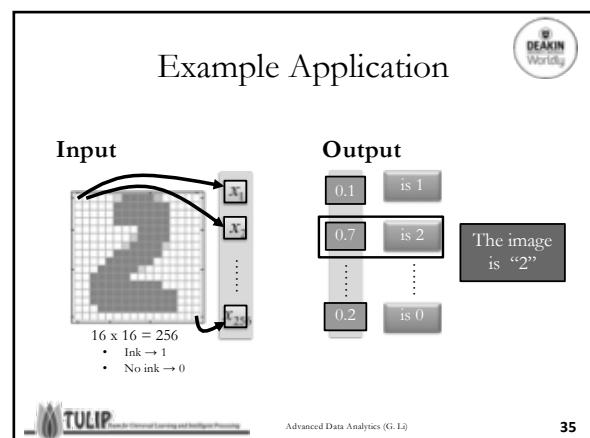
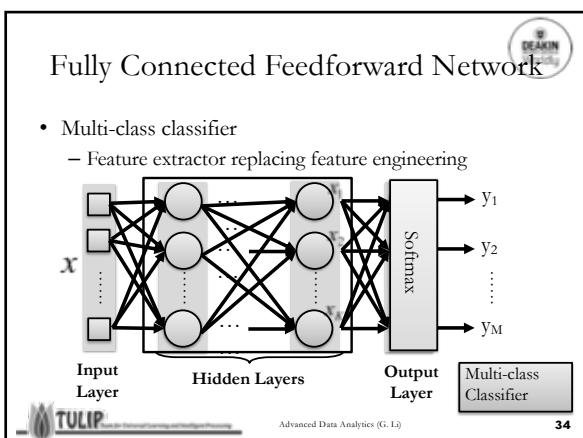
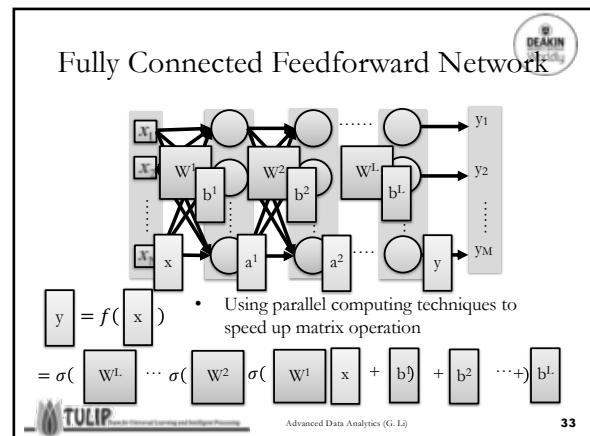
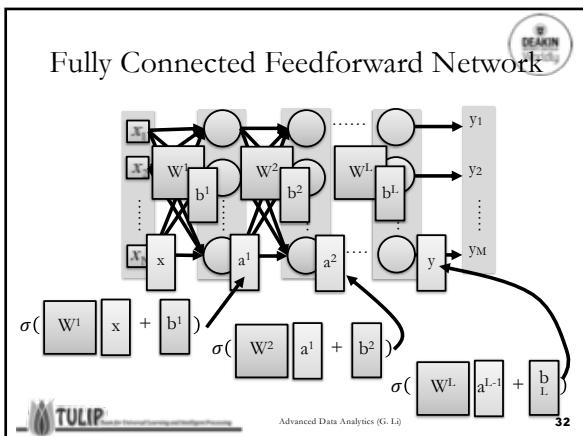
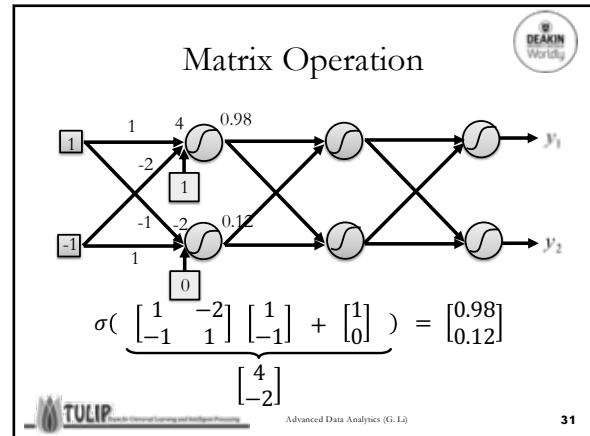
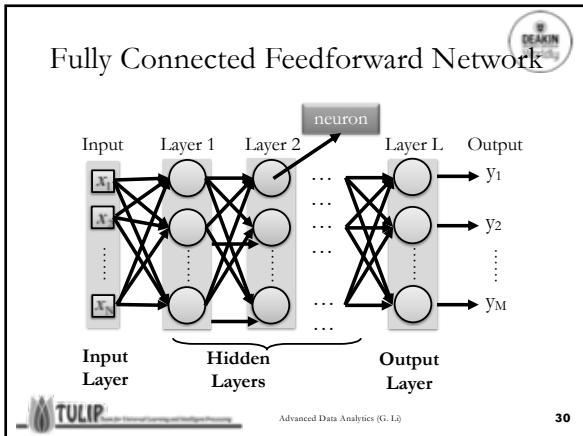
keynote of Rich Caruana at ASRU 2015

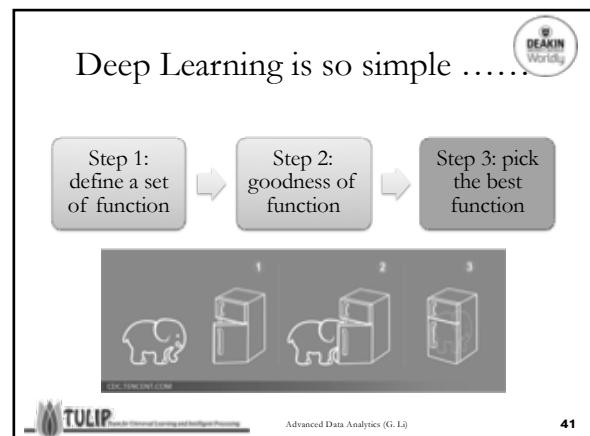
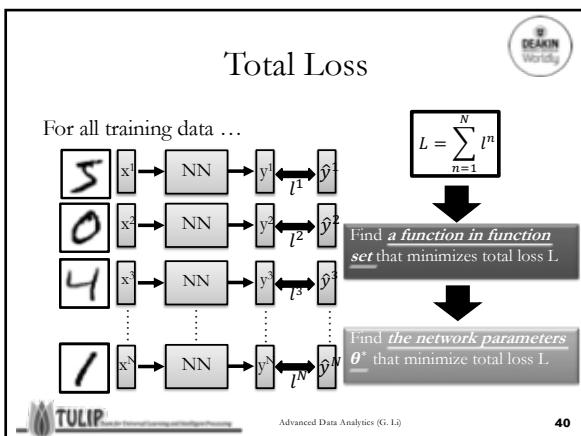
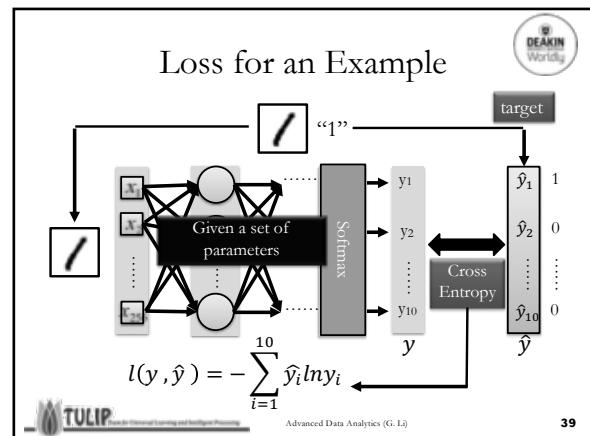
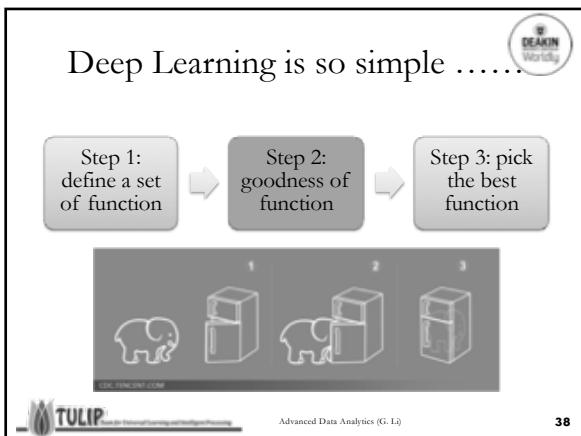
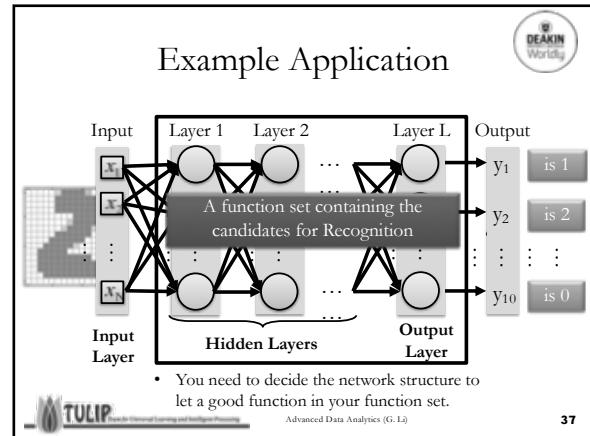
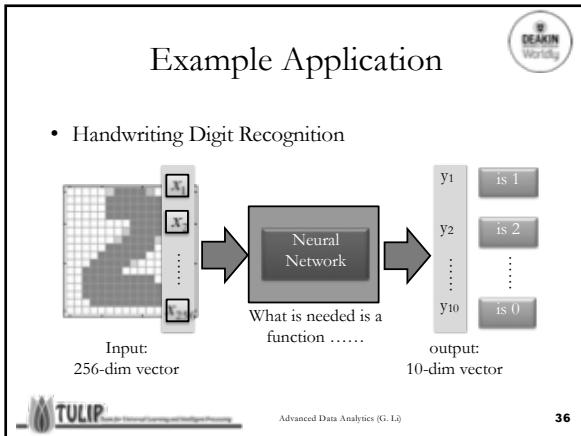
23

TULIP: Tools for Unsupervised Learning and Intelligent Processing

Advanced Data Analytics (G. Li)

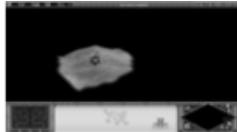




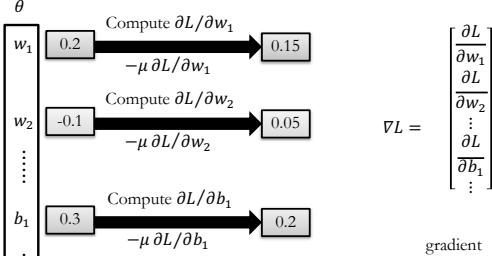


Gradient Descent

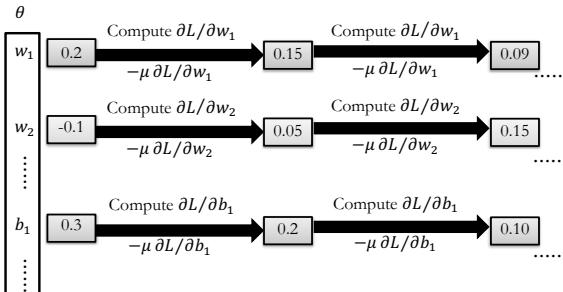
- Gradient Descent is the “learning” of machines in deep learning
- Even alpha go using this approach.



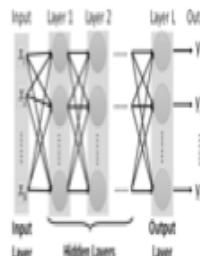
Gradient Descent



Gradient Descent



Back Propagation



Gradient Descent

Network parameters $\theta = \{w_1, w_2, \dots, b_1, b_2, \dots\}$

Starting Parameters $\theta^0 \rightarrow \theta^1 \rightarrow \theta^2 \rightarrow \dots$

$\nabla L(\theta)$

$$\begin{bmatrix} \frac{\partial L(\theta)}{\partial w_1} \\ \frac{\partial L(\theta)}{\partial w_2} \\ \vdots \\ \frac{\partial L(\theta)}{\partial b_1} \\ \frac{\partial L(\theta)}{\partial b_2} \\ \vdots \end{bmatrix}$$

Compute $\nabla L(\theta^0)$ $\theta^1 = \theta^0 - \eta \nabla L(\theta^0)$

Compute $\nabla L(\theta^1)$ $\theta^2 = \theta^1 - \eta \nabla L(\theta^1)$

Millions of parameters

To compute the gradients efficiently, we use **backpropagation**.

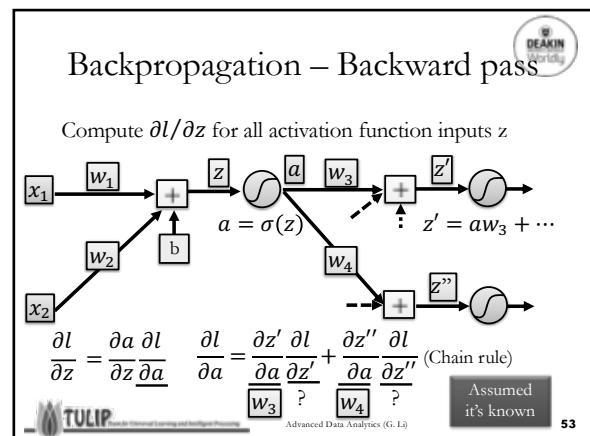
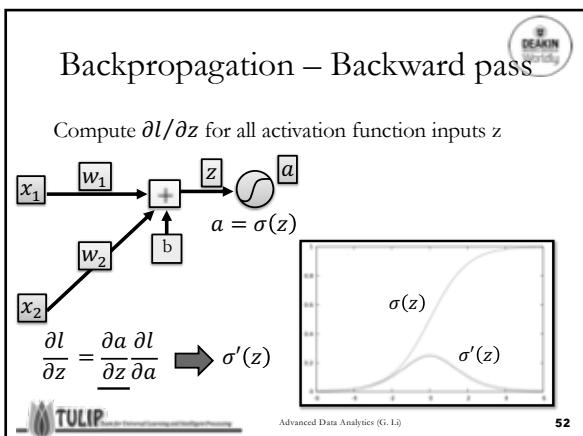
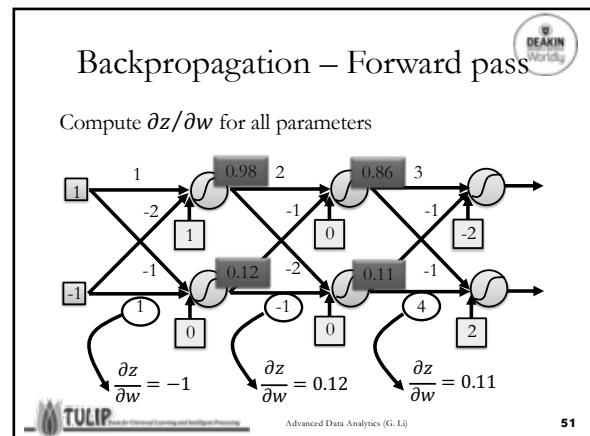
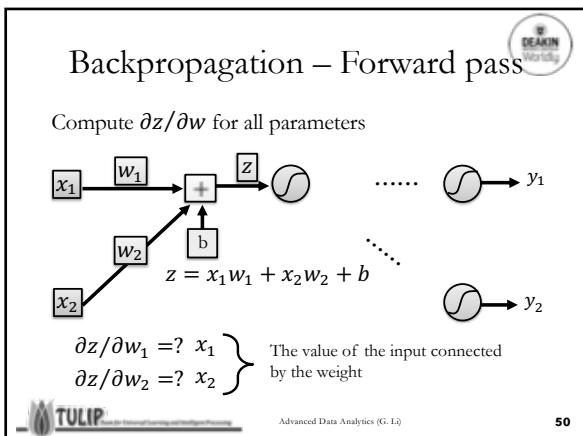
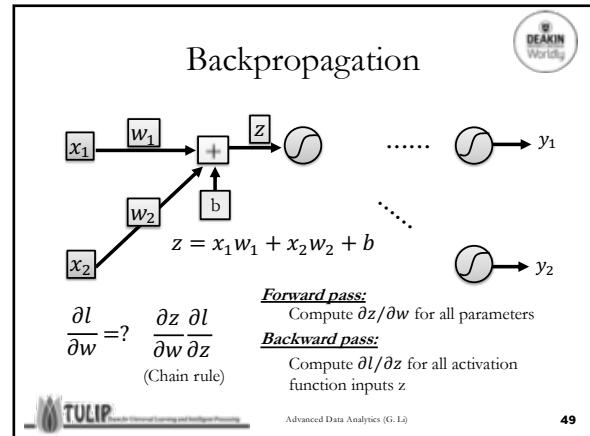
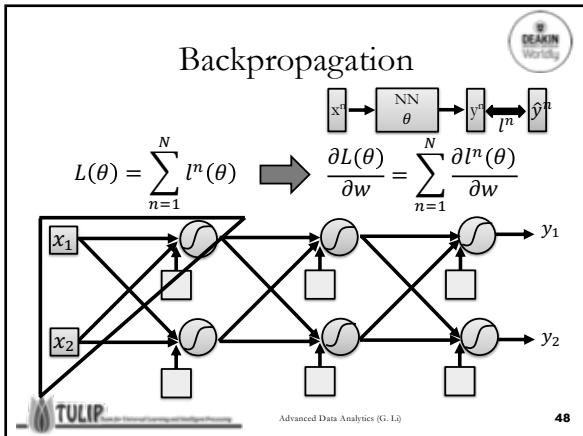
Chain Rule

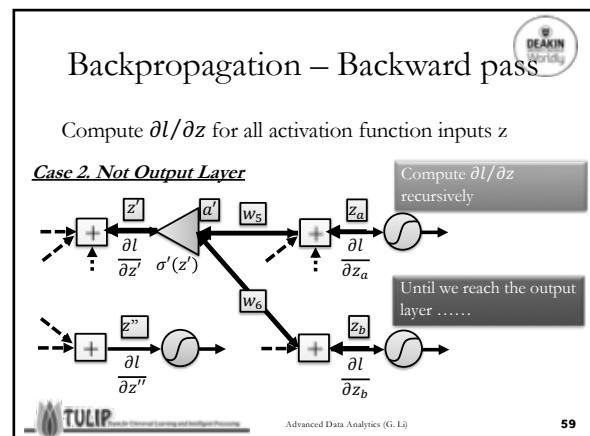
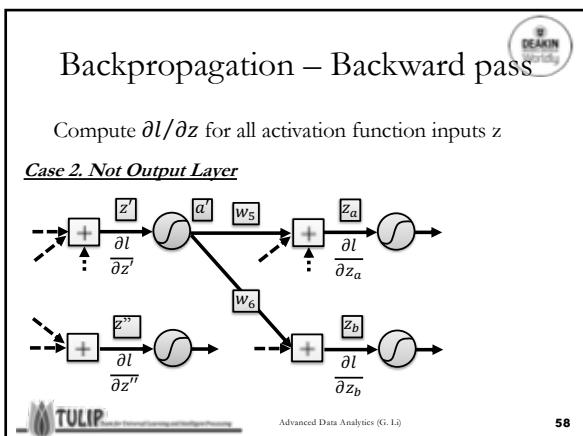
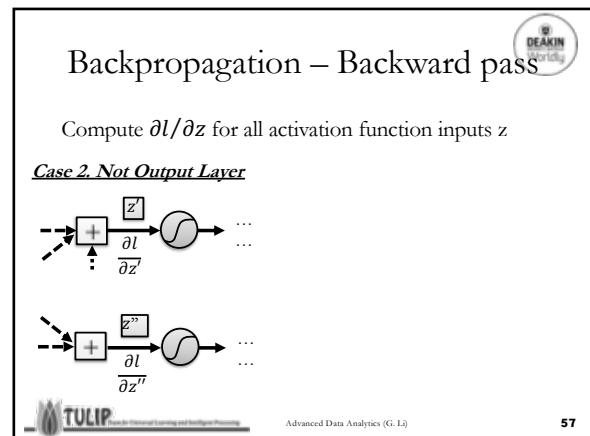
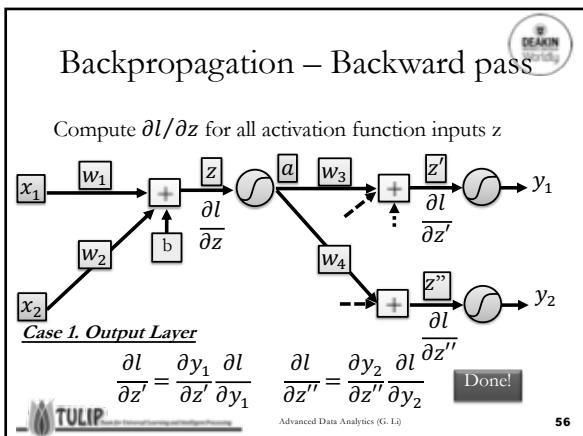
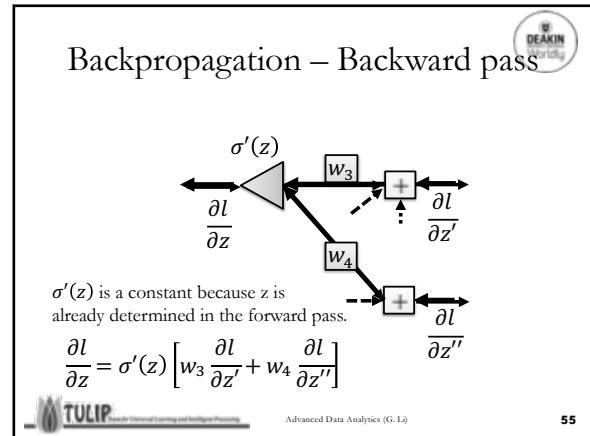
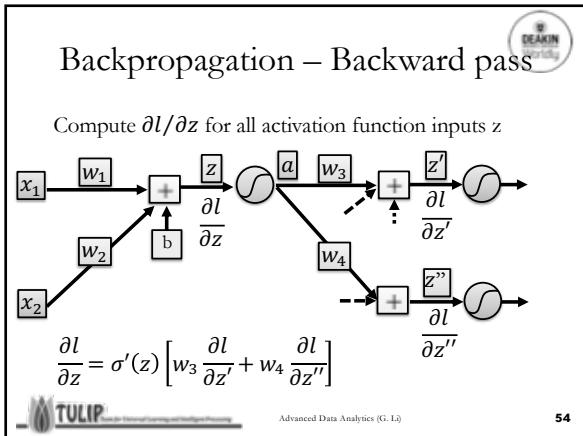
Case 1 $y = g(x) \quad z = h(y)$

$$\Delta x \rightarrow \Delta y \rightarrow \Delta z \quad \frac{dz}{dx} = \frac{dz}{dy} \frac{dy}{dx}$$

Case 2 $x = g(s) \quad y = h(s) \quad z = k(x, y)$

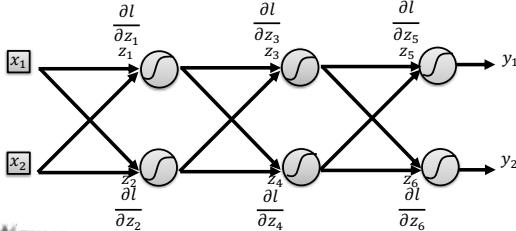
$$\Delta s \xrightarrow{\Delta x} \Delta z \quad \frac{dz}{ds} = \frac{\partial z}{\partial x} \frac{dx}{ds} + \frac{\partial z}{\partial y} \frac{dy}{ds}$$





Backpropagation – Backward Pass

Compute $\partial l / \partial z$ for all activation function inputs z
 Compute $\partial l / \partial z$ from the output layer

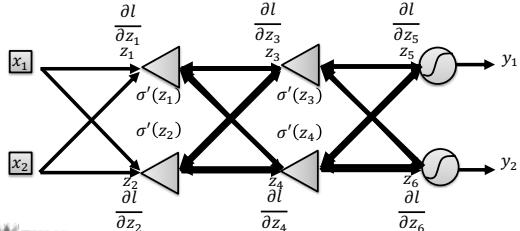


Advanced Data Analytics (G. Li)

60

Backpropagation – Backward Pass

Compute $\partial l / \partial z$ for all activation function inputs z
 Compute $\partial l / \partial z$ from the output layer

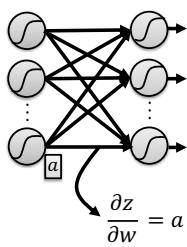


Advanced Data Analytics (G. Li)

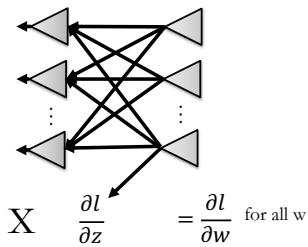
61

Backpropagation – Summary

Forward Pass



Backward Pass



Advanced Data Analytics (G. Li)

62

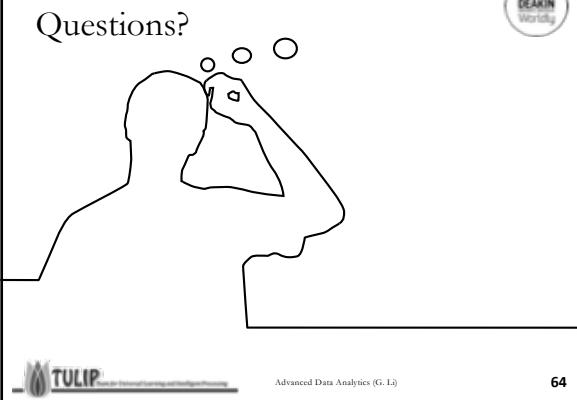
This Week's Readings

- Deep Learning: Theoretical Motivations (Yoshua Bengio)
 - http://videolectures.net/deeplearning2015_bengio_theoretical_motivations/
- Connections between physics and deep learning
 - <https://www.youtube.com/watch?v=5MdSE-N0bxw>
- Why Deep Learning Works: Perspectives from Theoretical Chemistry
 - <https://www.youtube.com/watch?v=kIbKHIPbxiU>



63

Questions?



Advanced Data Analytics (G. Li)

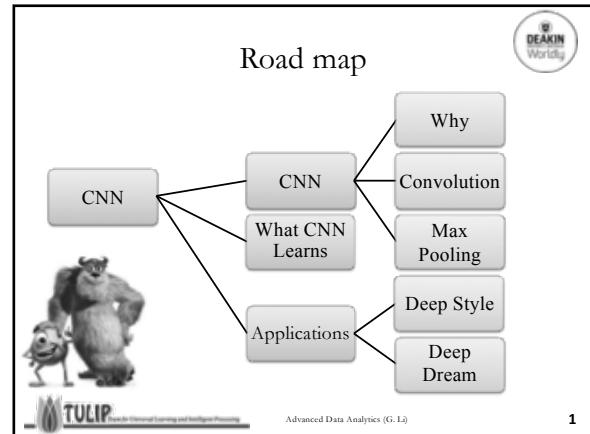
64

Lecture Notes on
Advanced Data Analytics

Module 05: Convolution Neural Network

Gang Li
School of Information Technology
Deakin University, VIC 3125, Australia





Convolution Neural Network



- Why CNN?
- Convolution
- Max Pooling

Advanced Data Analytics (G. Li) 2

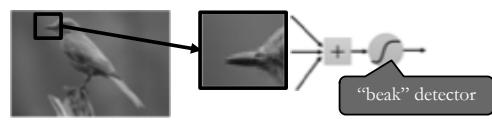


Why CNN for Image (1)

- Some patterns are much smaller than the whole image

A neuron does not have to see the whole image to discover the pattern.

Connecting to small region with less parameters



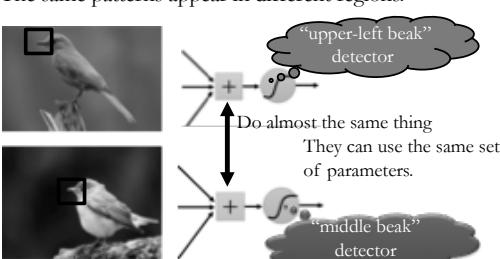
“beak” detector

Advanced Data Analytics (G. Li) 3




Why CNN for Image (2)

- The same patterns appear in different regions.



Do almost the same thing
They can use the same set of parameters.

“upper-left beak” detector

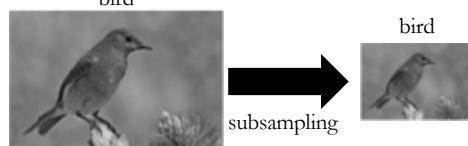
“middle beak” detector

Advanced Data Analytics (G. Li) 4



Why CNN for Image (3)

- Subsampling the pixels will not change the object



bird

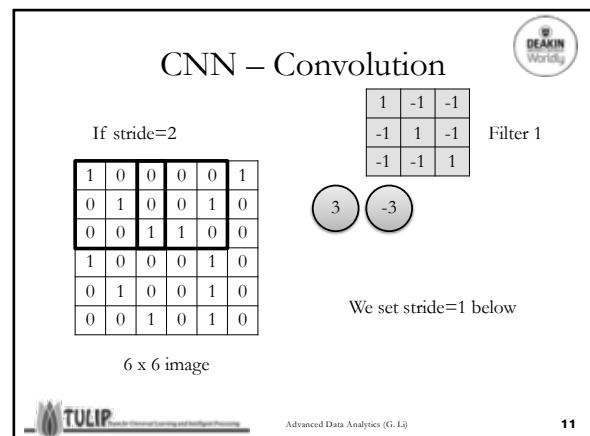
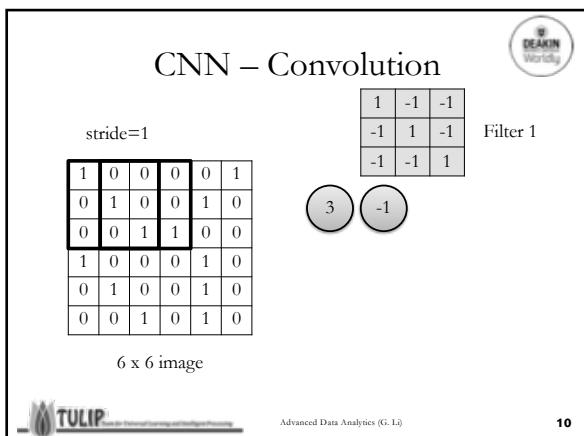
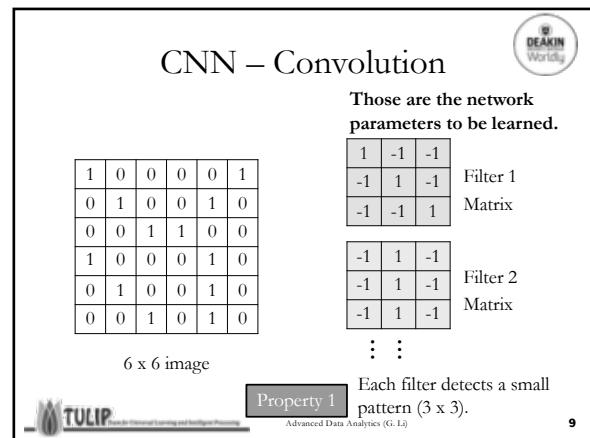
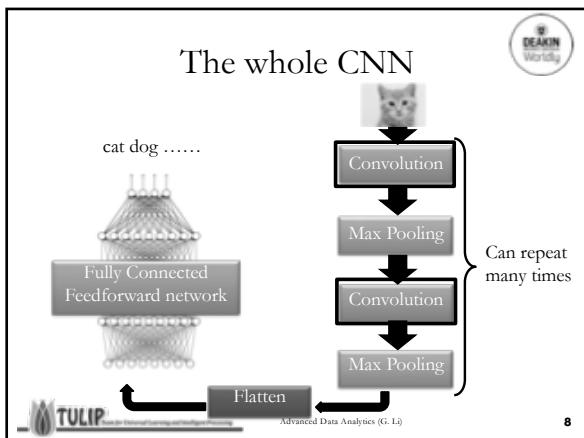
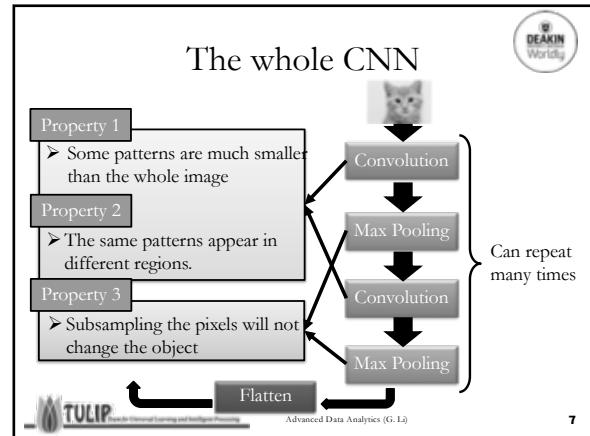
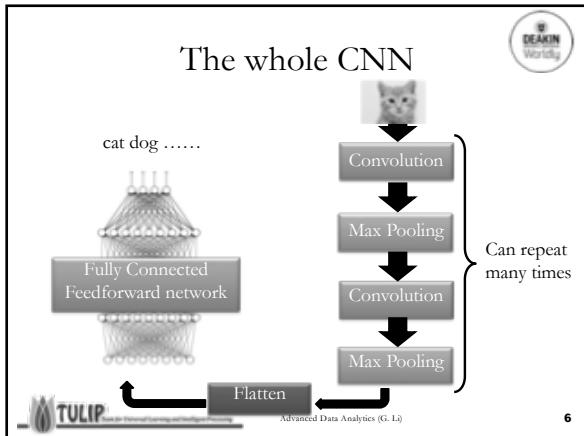
subsampling

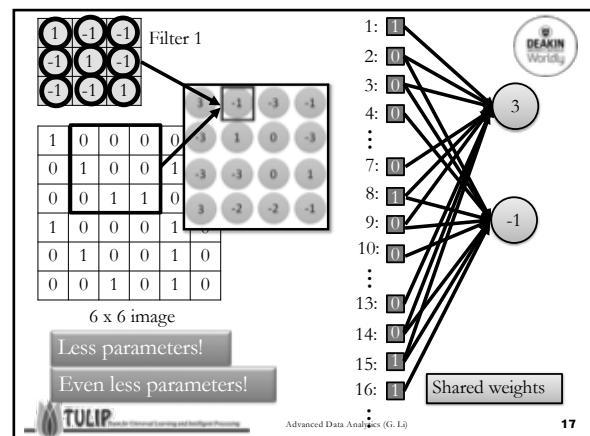
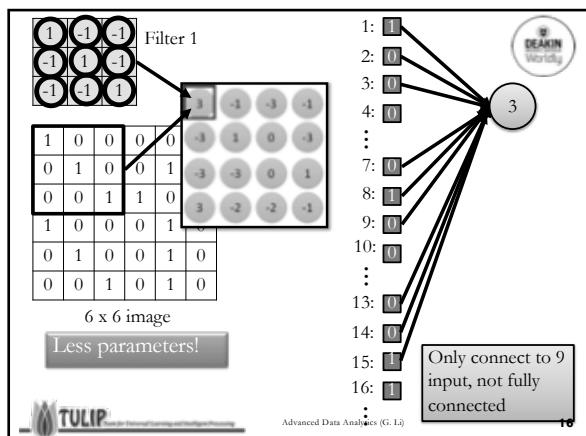
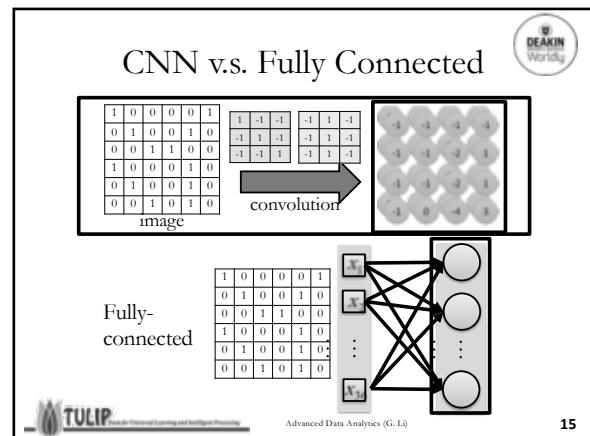
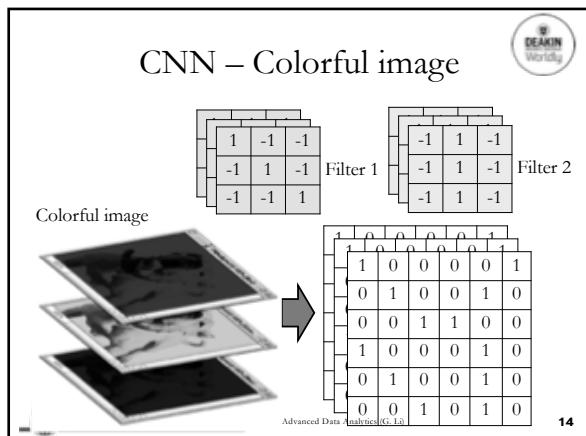
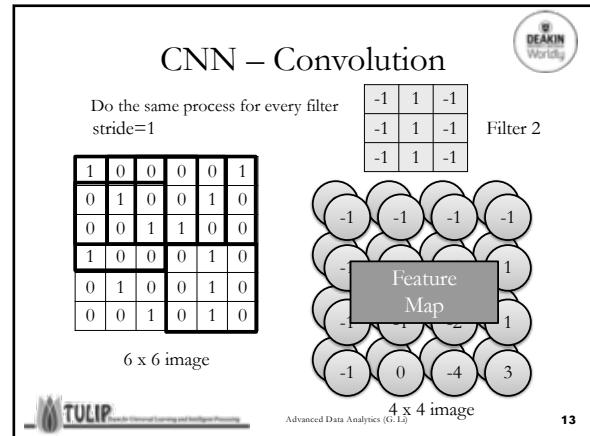
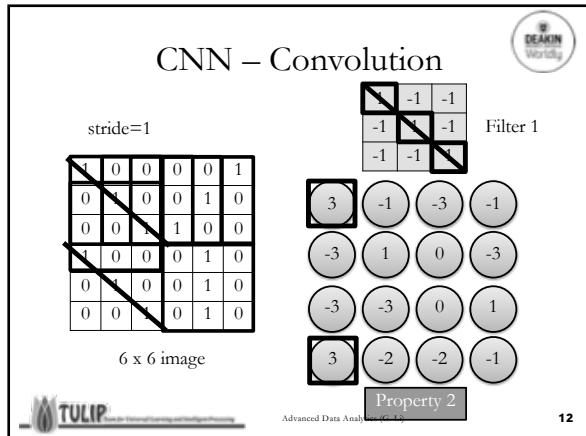
bird

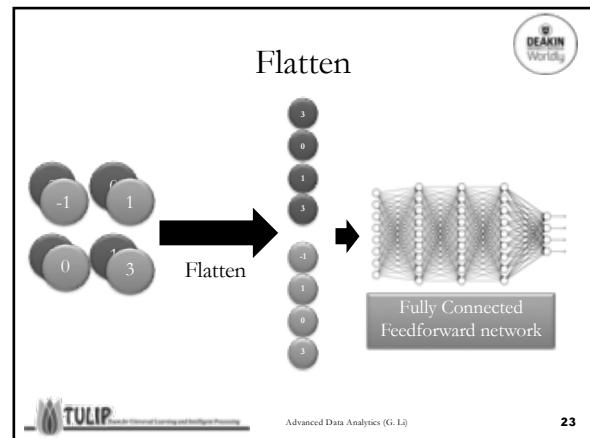
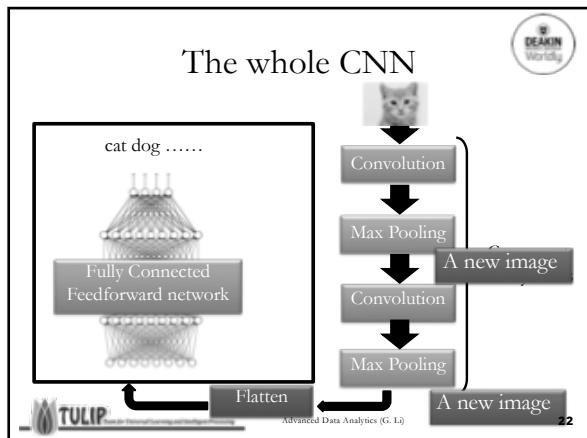
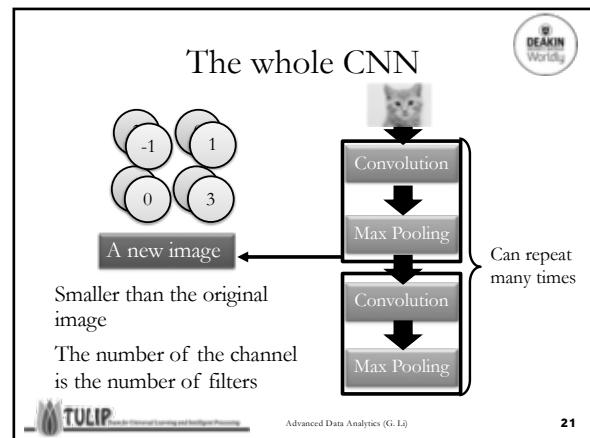
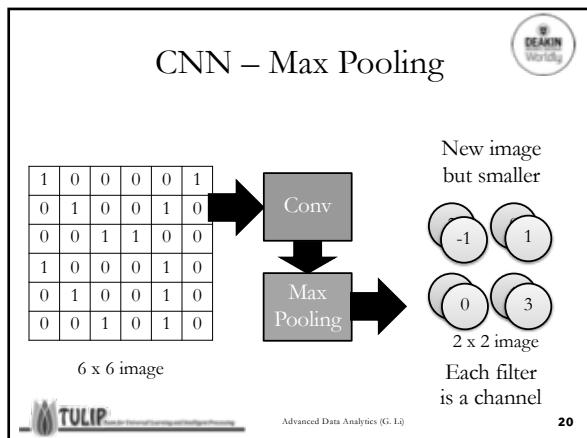
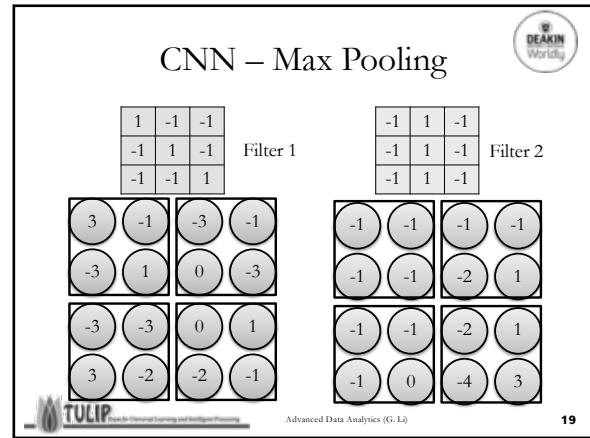
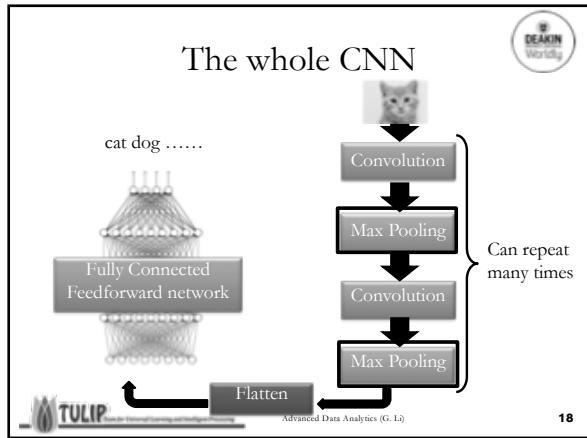
We can subsample the pixels to make image smaller
→ Less parameters for the network to process the image

Advanced Data Analytics (G. Li) 5







What CNN Learns?

- What CNN Learns?

Advanced Data Analytics (G. Li) 24

First Convolution Layer

- Typical-looking filters on the trained first layer

11 x 11 (AlexNet)

<http://cs231n.github.io/understanding-cnn/>

Advanced Data Analytics (G. Li) 25

How about higher layers?

- Which images make a specific neuron activate

Ross Girshick, Jeff Donahue, Trevor Darrell, Jitendra Malik, "Rich feature hierarchies for accurate object detection and semantic segmentation", CVPR, 2014

Advanced Data Analytics (G. Li) 26

What does CNN learn?

The output of the k-th filter is a 11 x 11 matrix.

Degree of the activation of the k-th filter:

$$a^k = \sum_{i=1}^{11} \sum_{j=1}^{11} a_{ij}^k$$

(gradient ascent)

$$x^* = \arg \max_x a^k$$

Advanced Data Analytics (G. Li) 27

What does CNN learn?

The output of the k-th filter is a 11 x 11 matrix.

Degree of the activation of the k-th filter:

$$a^k = \sum_{i=1}^{11} \sum_{j=1}^{11} a_{ij}^k$$

(gradient ascent)

$$x^* = \arg \max_x a^k$$

Advanced Data Analytics (G. Li) 28

What does CNN learn?

Find an image maximizing the output of neuron:

$$x^* = \arg \max_x a^j$$

Each figure corresponds to a neuron

Advanced Data Analytics (G. Li) 29

What does CNN learn?

$$x^* = \arg \max_x y^i$$

$$x^* = \arg \max_x \left(y^i - \sum_{i,j} |x_{ij}| \right)$$

Advanced Data Analytics (G. Li) 30

What does CNN learn?

- Karen Simonyan, Andrea Vedaldi, Andrew Zisserman, "Deep Inside Convolutional Networks: Visualising Image Classification Models and Saliency Maps", ICLR, 2014

Advanced Data Analytics (G. Li) 31

What does CNN learn?

$|\frac{\partial y_k}{\partial x_{ij}}|$

y_k : the predicted class of the model

Pixel x_{ij}

- Karen Simonyan, Andrea Vedaldi, Andrew Zisserman, "Deep Inside Convolutional Networks: Visualising Image Classification Models and Saliency Maps", ICLR, 2014

Advanced Data Analytics (G. Li) 32

What does CNN learn?

- Reference: Zeiler, M. D., & Fergus, R. (2014). Visualizing and understanding convolutional networks. In Computer Vision–ECCV 2014 (pp. 818–833)

Advanced Data Analytics (G. Li) 33

Applications

- Deep Dream
- Deep Style

Advanced Data Analytics (G. Li) 34

Deep Dream Generator

<https://deepr dreamgenerator.com>

guillermoponce 12 hours ago

Advanced Data Analytics (G. Li) 35

Deep Style

- Given a photo, make its style like famous paintings



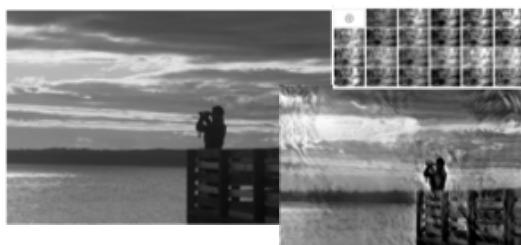
Advanced Data Analytics (G. Li)

36



Deep Style

- Given a photo, make its style like famous paintings



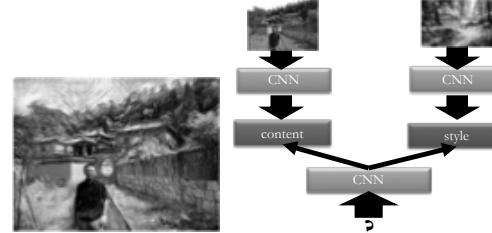
Advanced Data Analytics (G. Li)

38



Deep Style

- Given a photo, make its style like famous paintings

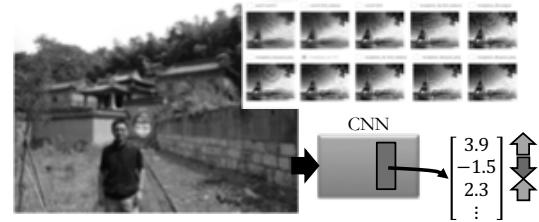
A Neural Algorithm of Artistic Style (<https://arxiv.org/abs/1508.06576>)

37



Deep Dream

- Given a photo, machine adds what it sees



• CNN exaggerates what it sees



Advanced Data Analytics (G. Li)

39



Deep Dream

- Given a photo, machine adds what it sees

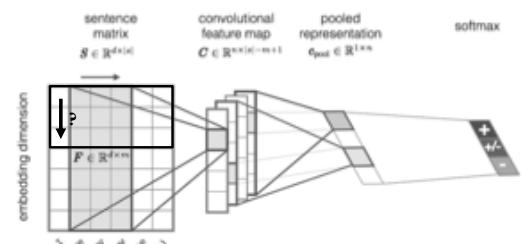


Advanced Data Analytics (G. Li)

40



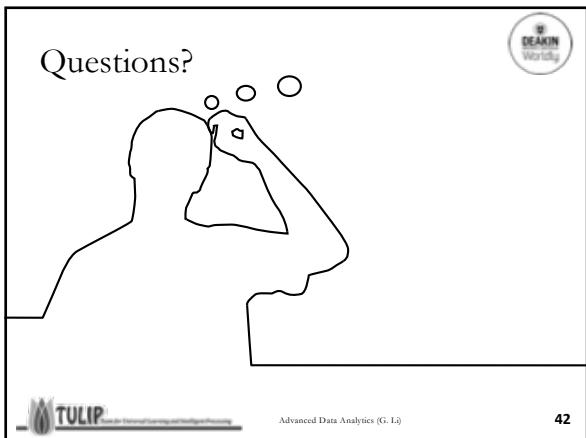
More Application: Text

Source of image:
<http://cagcens.csail.mit.edu/vision/doc/download?doi=10.1177/0363685887001001001.pdf>

Advanced Data Analytics (G. Li)

41





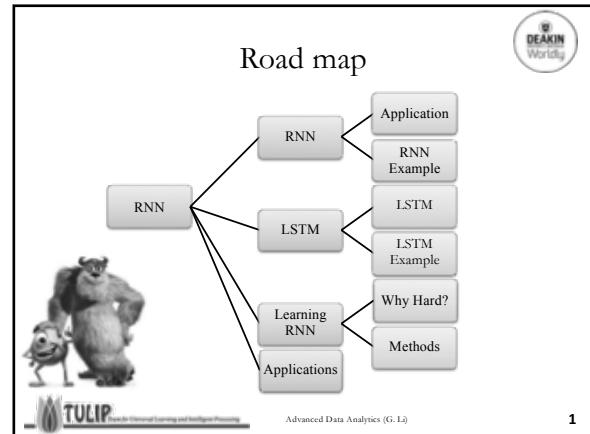
Lecture Notes on
Advanced Data Analytics

Module 05: Recurrent Neural Network (RNN)

Gang Li
School of Information Technology
Deakin University, VIC 3125, Australia



TULIP Tool for Universal Learning and Intelligent Processing



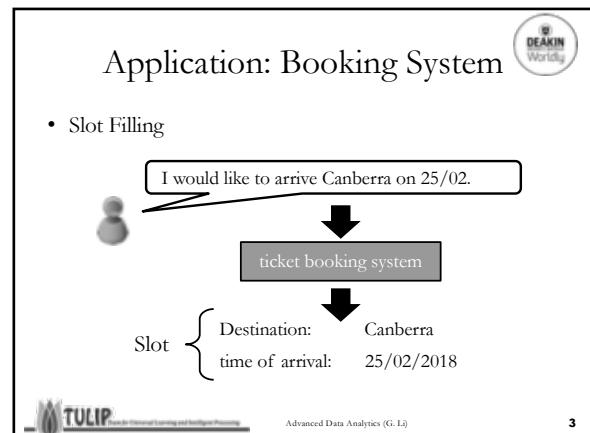
Recurrent Neural Network



- Why RNN?
- RNN
- RNN Example

TULIP Tool for Universal Learning and Intelligent Processing

Advanced Data Analytics (G. Li) 2



Application: Booking System

- Solving slot filling by FFN?
- Input:
 - Word represented as a vector

Canberra → [x₁, x₂] → y₁, y₂

TULIP Tool for Universal Learning and Intelligent Processing

Advanced Data Analytics (G. Li) 4

One Hot Encoding

- The vector is lexicon size.
- Each dimension corresponds to a word in the lexicon
- The dimension for the word is 1, and others are 0

lexicon = {apple, bag, cat, dog, elephant}

apple	=	[1 0 0 0 0]
bag	=	[0 1 0 0 0]
cat	=	[0 0 1 0 0]
dog	=	[0 0 0 1 0]
elephant	=	[0 0 0 0 1]

TULIP Tool for Universal Learning and Intelligent Processing

Advanced Data Analytics (G. Li) 5

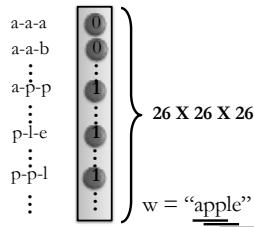
Beyond One Hot Encoding

Dimension for “Other”

apple	0
bag	0
cat	0
dog	0
elephant	0
“other”	1

w = “Gandalf” w = “Sauron”

Word hashing



Advanced Data Analytics (G. Li)

Application: Booking System

Solving slot filling by FFN?

– Input:

- Word represented as a vector

– Output:

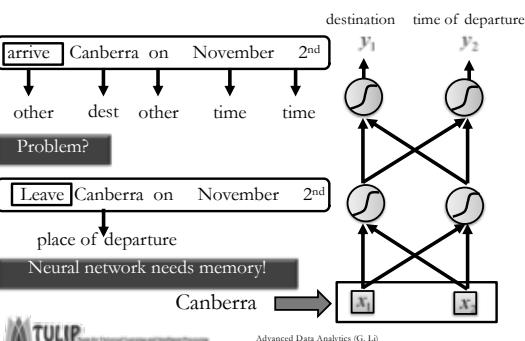
- Probability distribution that the input word belonging to the slots

Canberra

Advanced Data Analytics (G. Li)

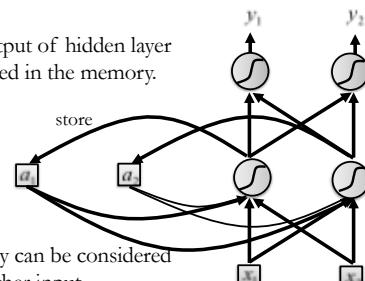
7

Example Application



Recurrent Neural Network (RNN)

The output of hidden layer are stored in the memory.



Advanced Data Analytics (G. Li)

9

Example

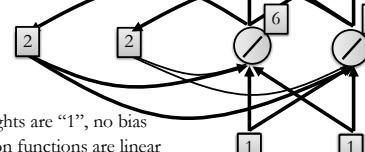
Input sequence: $\begin{bmatrix} 1 \\ 1 \end{bmatrix} \begin{bmatrix} 1 \\ 1 \end{bmatrix} \begin{bmatrix} 2 \\ 2 \end{bmatrix} \dots \dots$ Output sequence: $\begin{bmatrix} 4 \\ 4 \end{bmatrix}$ Initial values $\begin{bmatrix} 0 \\ 0 \end{bmatrix}$ All the weights are “1”, no bias
All activation functions are linear

TULIP

Advanced Data Analytics (G. Li)

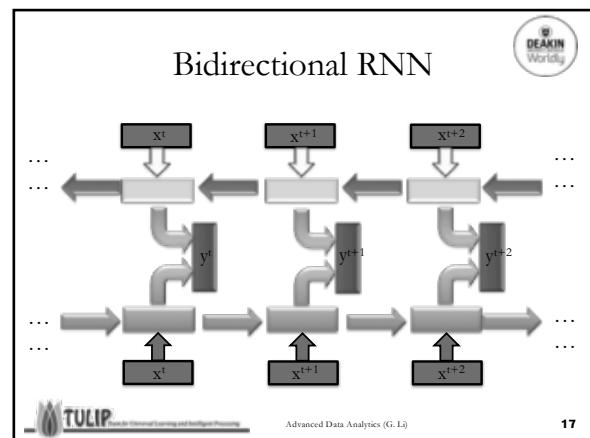
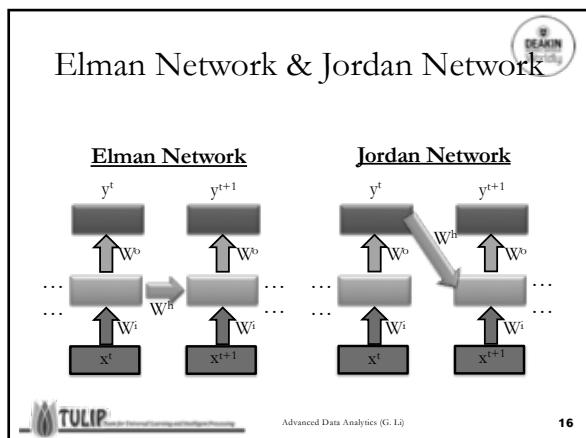
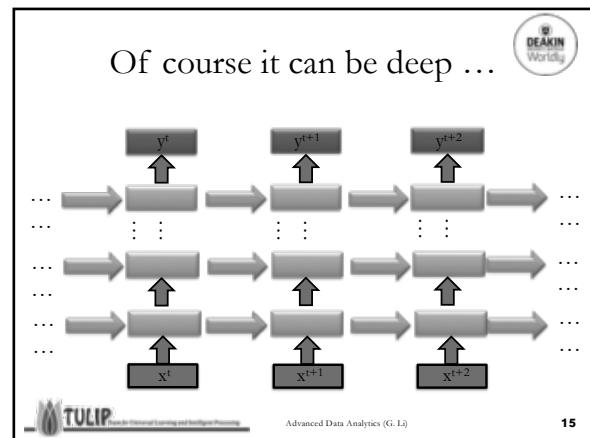
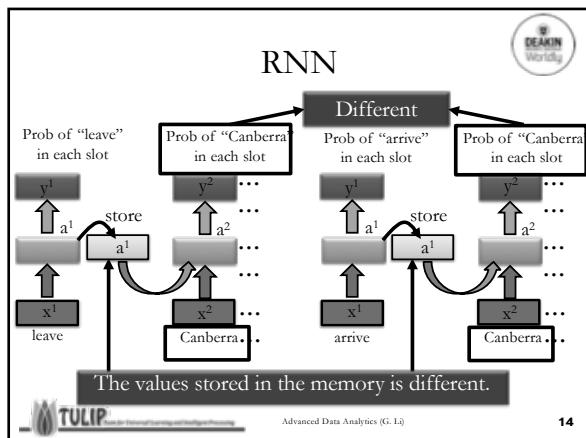
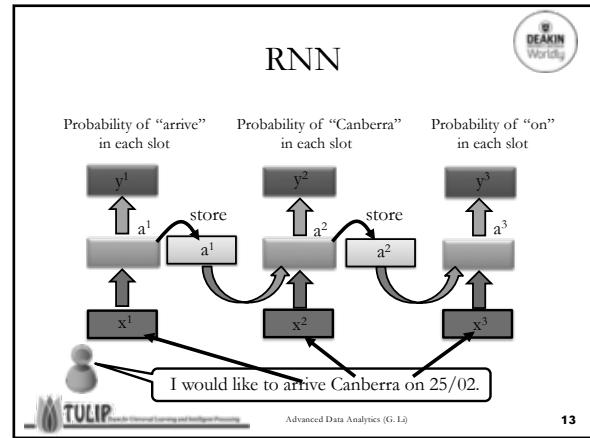
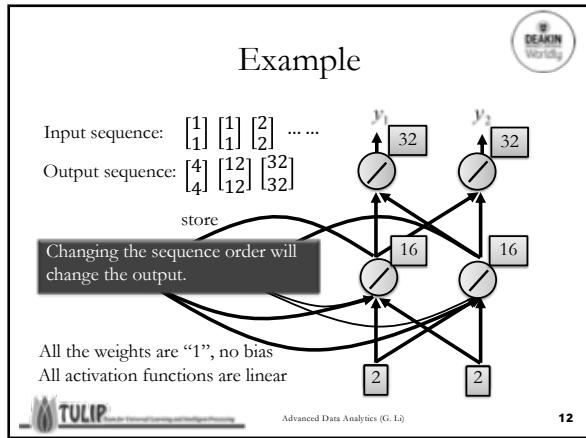
10

Example

Input sequence: $\begin{bmatrix} 1 \\ 1 \end{bmatrix} \begin{bmatrix} 1 \\ 1 \end{bmatrix} \begin{bmatrix} 2 \\ 2 \end{bmatrix} \dots \dots$ Output sequence: $\begin{bmatrix} 4 \\ 4 \end{bmatrix} \begin{bmatrix} 12 \\ 12 \end{bmatrix}$ 

Advanced Data Analytics (G. Li)

11

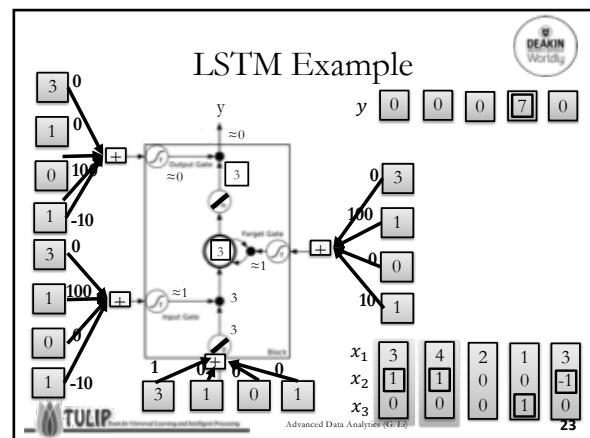
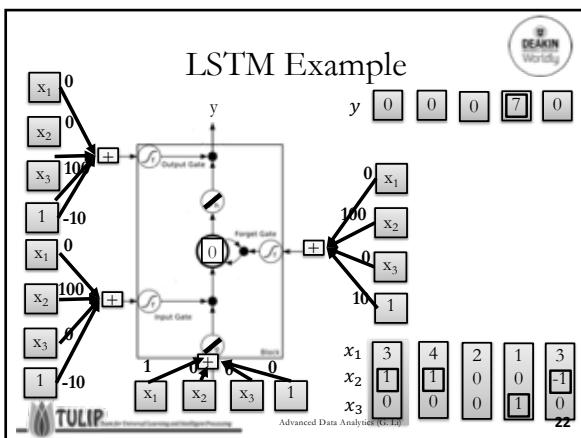
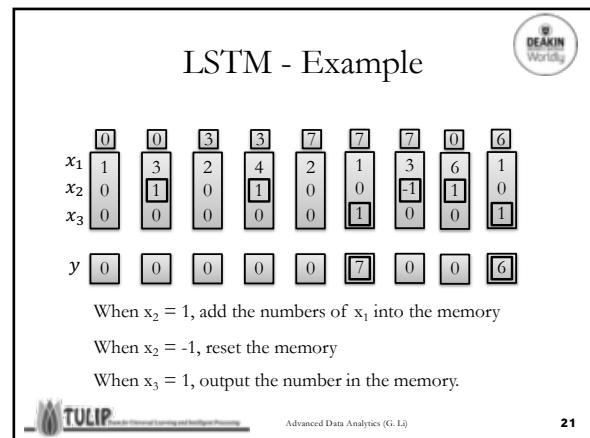
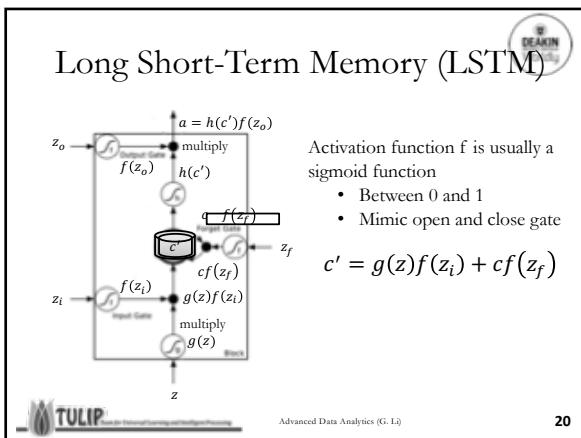
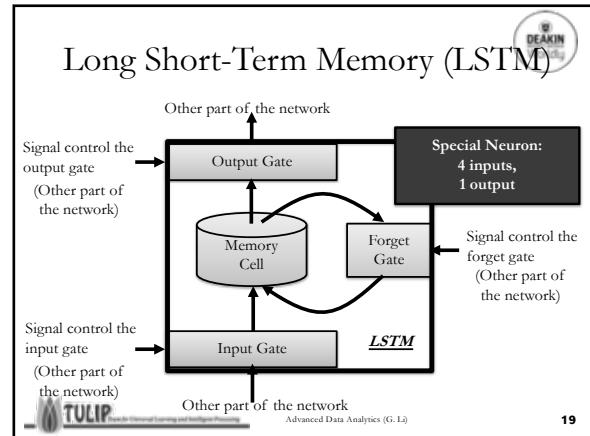


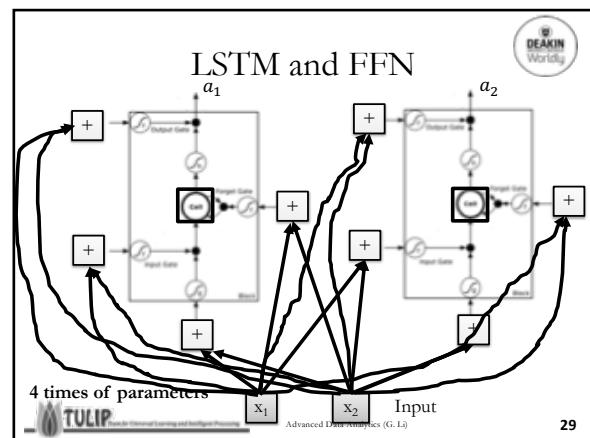
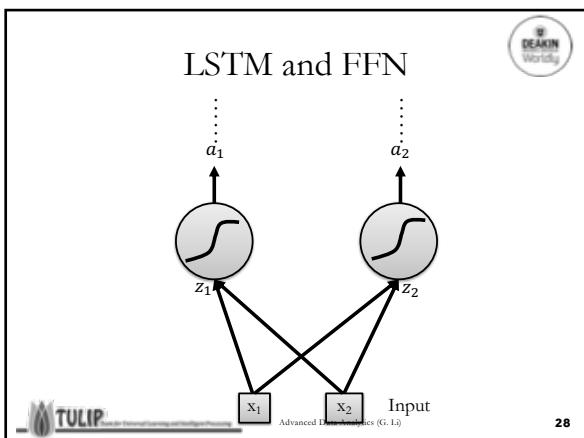
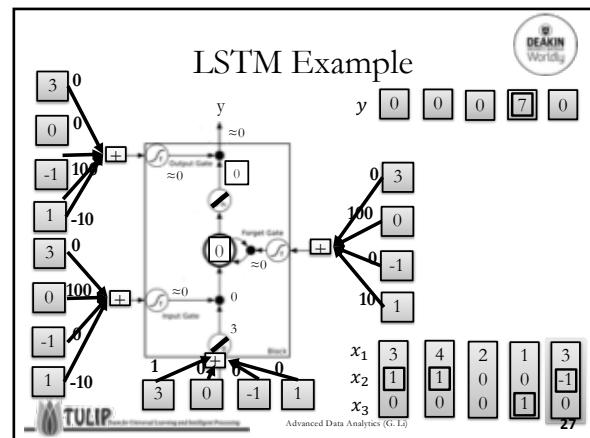
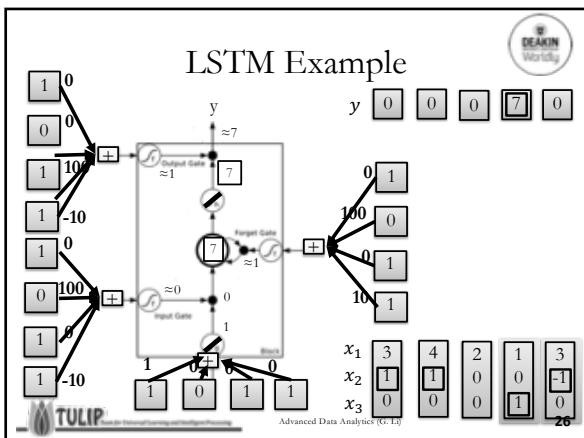
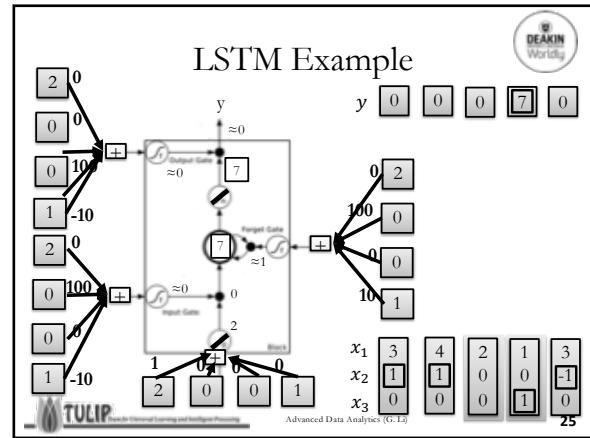
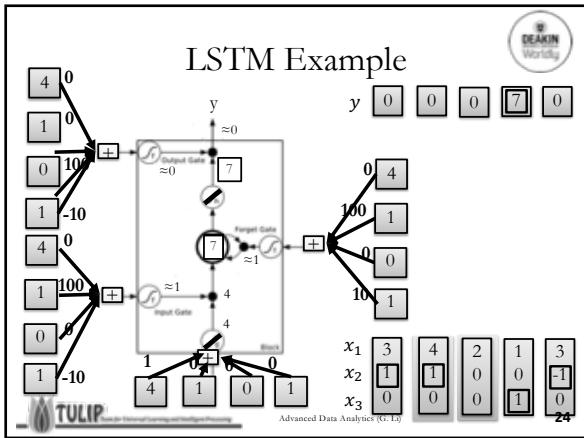
LSTM

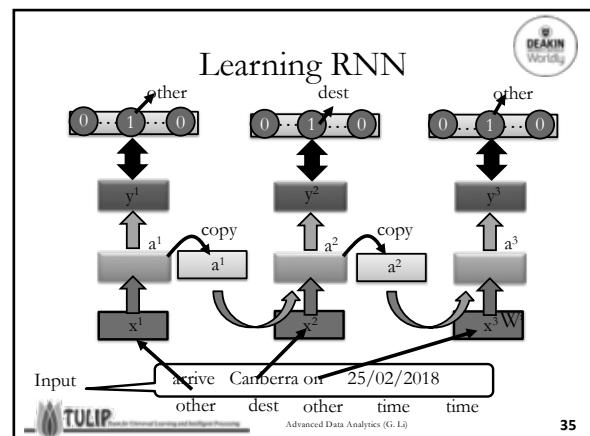
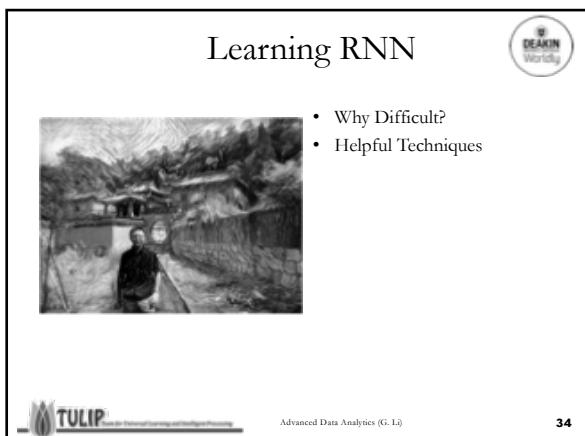
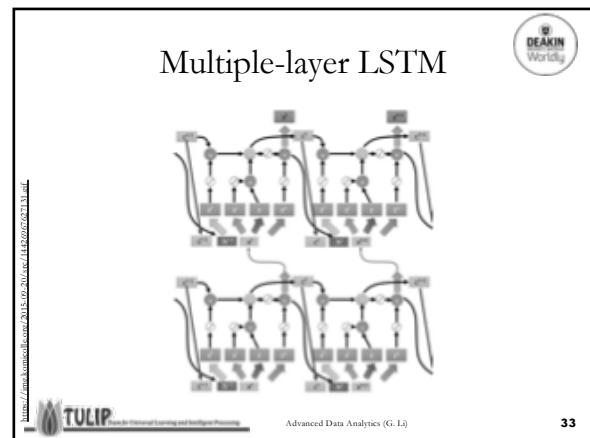
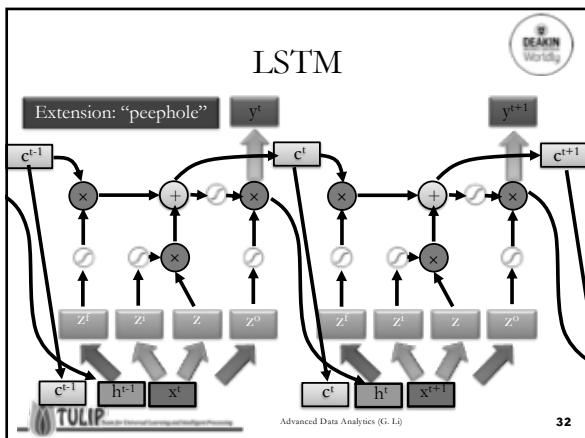
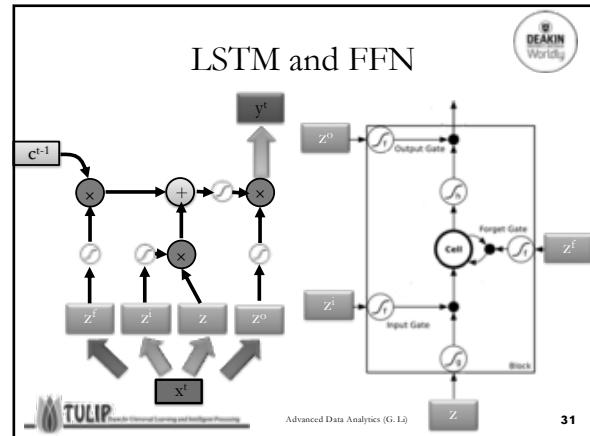
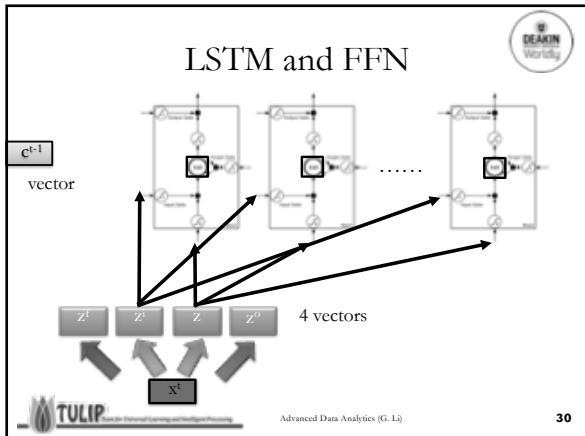
- Why LSTM?
- LSTM
- LSTM and FFN

Advanced Data Analytics (G. Li)

18







Learning RNN

Backpropagation through time (BPTT)

$$w \leftarrow w - \eta \frac{\partial L}{\partial w}$$

TULIP: Tools for Unsupervised Learning and Intelligent Processing Advanced Data Analytics (G. Li) 36

Learning RNN

- RNN-based network is not always easy to learn

Real experiments on Language modeling

TULIP: Tools for Unsupervised Learning and Intelligent Processing Advanced Data Analytics (G. Li) 37

Learning RNN

The error surface is either very flat or very steep.

[Razvan Pascanu, ICML'13]

TULIP: Tools for Unsupervised Learning and Intelligent Processing Advanced Data Analytics (G. Li) 38

Why Difficult?

$w = 1 \rightarrow y^{1000} = 1$	$w = 1.01 \rightarrow y^{1000} \approx 20000$	Large $\frac{\partial L}{\partial w}$	Small Learning rate?
$w = 0.99 \rightarrow y^{1000} \approx 0$	$w = 0.01 \rightarrow y^{1000} \approx 0$	small $\frac{\partial L}{\partial w}$	Large Learning rate?

Toy Example

TULIP: Tools for Unsupervised Learning and Intelligent Processing Advanced Data Analytics (G. Li) 39

Helpful Techniques

- Long Short-term Memory (LSTM)
 - Can deal with gradient vanishing (not gradient explode)
 - Memory and input are added
 - The influence never disappears unless forget gate is closed

No Gradient vanishing
(If forget gate is opened.)

TULIP: Tools for Unsupervised Learning and Intelligent Processing Advanced Data Analytics (G. Li) 40

Helpful Techniques

Clockwise RNN

[Jan Koutnuk, JMLR'14]

Structurally Constrained Recurrent Network (SCRN)

[Tomas Mikolov, ICLR'15]

Vanilla RNN Initialized with Identity matrix + ReLU activation function [Quoc V. Le, arXiv'15]

- Outperform or be comparable with LSTM in 4 different tasks

TULIP: Tools for Unsupervised Learning and Intelligent Processing Advanced Data Analytics (G. Li) 41

RNN/LSTM Applications



- Applications
- Attention-based Model
- Visual Question Answering

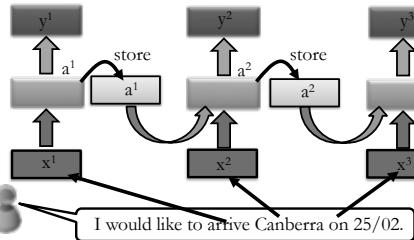


Advanced Data Analytics (G. Li)

42

Applications: Many to Many

Probability of “arrive” in each slot Probability of “Canberra” in each slot Probability of “on” in each slot



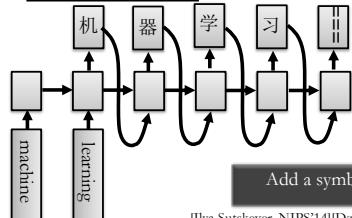
Advanced Data Analytics (G. Li)

43

Many to Many (No Limitation)

- Both input and output are both sequences with different lengths → Sequence to sequence learning

Machine Translation



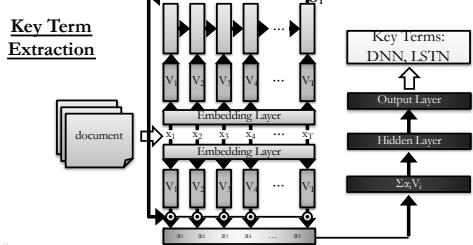
[Ilya Sutskever, NIPS'14][Dzmitry Bahdanau, arXiv'15]

Advanced Data Analytics (G. Li)

44

Applications: Many to One

- Input is a vector sequence, but output is only one vector

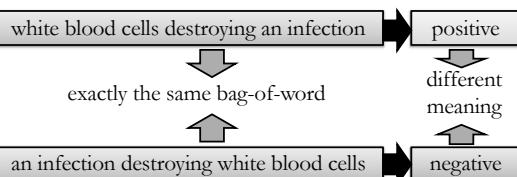


Advanced Data Analytics (G. Li)

45

Sequence-to-sequence Auto-encoder - Text

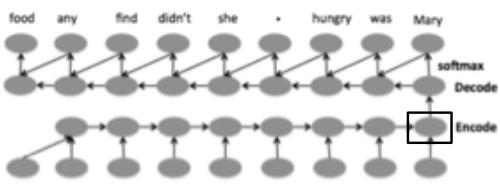
- To understand the meaning of a word sequence, the order of the words can not be ignored.



Advanced Data Analytics (G. Li)

46

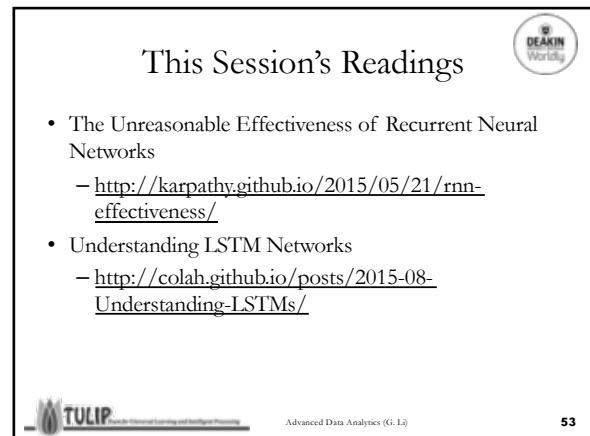
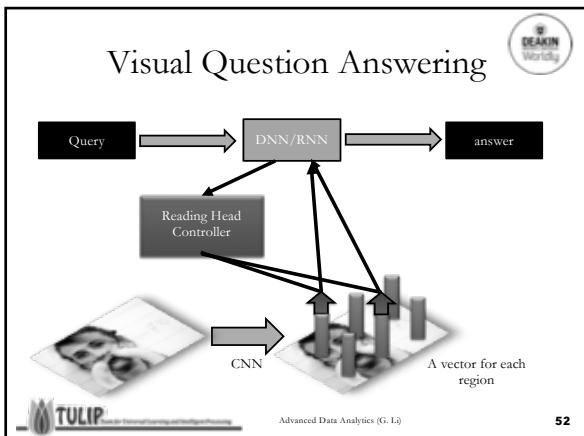
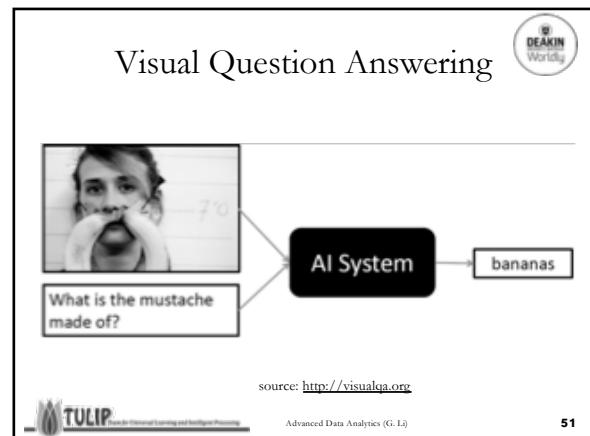
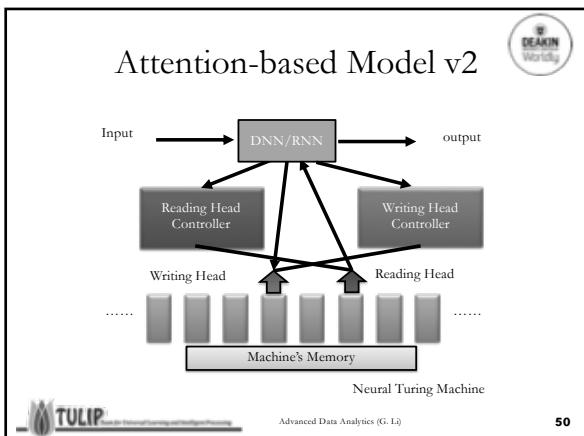
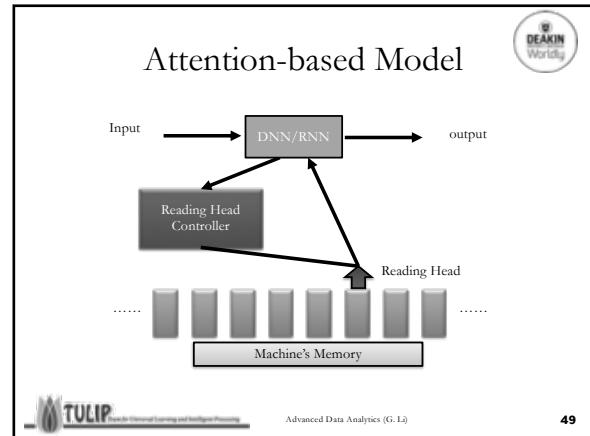
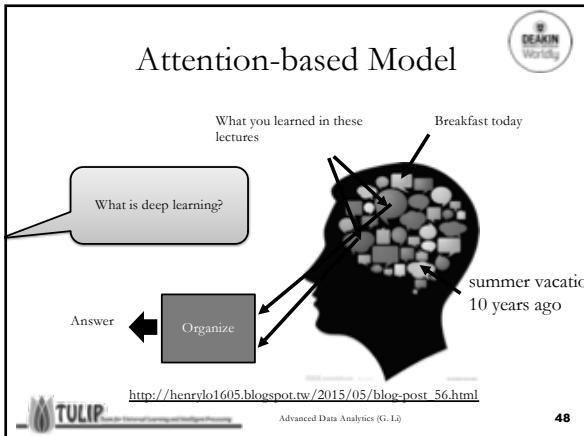
Sequence-to-sequence Auto-encoder - Text

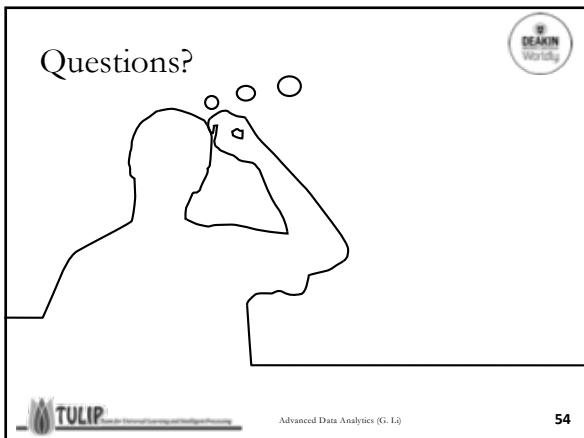


- Li, Jiwei, Minh-Thang Luong, and Dan Jurafsky. "A hierarchical neural autoencoder for paragraphs and documents." *arXiv preprint arXiv:1506.01057*(2015).

Advanced Data Analytics (G. Li)

47



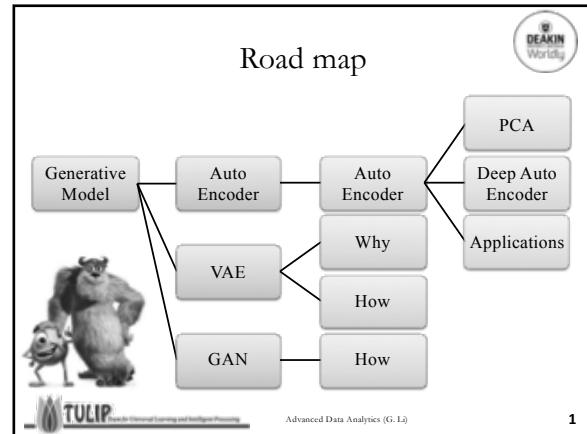


Lecture Notes on
Advanced Data Analytics

Module 05: Generative Model

Gang Li
School of Information Technology
Deakin University, VIC 3125, Australia

TULIP Team for Unsupervised Learning and Intelligence Processing



Generative Models

- Unsupervised Learning
- Generation

TULIP Team for Unsupervised Learning and Intelligence Processing

Unsupervised Learning

- “We expect unsupervised learning to become far more important in the longer term. Human and animal learning is largely unsupervised: we discover the structure of the world by observing it, not by being told the name of every object.”
LeCun, Bengio, Hinton, Nature 2015
- As I've said in previous statements: most of human and animal learning is unsupervised learning. If intelligence was a cake, unsupervised learning would be the cake, supervised learning would be the icing on the cake, and reinforcement learning would be the cherry on the cake. We know how to make the icing and the cherry, but we don't know how to make the cake.
• Yann LeCun, March 14, 2016 (Facebook)

TULIP Team for Unsupervised Learning and Intelligence Processing

Unsupervised Learning

Dimension Reduction

only having function input

Generation

only having function output

Random numbers

TULIP Team for Unsupervised Learning and Intelligence Processing

Generation

- Generative Models:
– <https://openai.com/blog/generative-models/>

What I cannot create,
I do not understand.

Richard Feynman

<https://www.youtube.com/What-did-Richard-Feynman-mean-when-he-said-What-I-cannot-create-I-do-not-understand>

TULIP Team for Unsupervised Learning and Intelligence Processing

Generation

Now v.s.
In the future
Machine draws a cat

TULIP Tool for Unsupervised Learning and Intelligence Processing

Advanced Data Analytics (G. Li)

6

Component-by-component

- Image generation

E.g. 3 x 3 images

Can be trained just with a large collection of images without any annotation

TULIP Tool for Unsupervised Learning and Intelligence Processing

Advanced Data Analytics (G. Li)

7

Component-by-component

- Image generation

E.g. 3 x 3 images

Can be trained just with a large collection of images without any annotation

TULIP Tool for Unsupervised Learning and Intelligence Processing

Advanced Data Analytics (G. Li)

8

PixelRNN

Real World

Aaron van den Oord, Nal Kalchbrenner, Koray Kavukcuoglu, Pixel Recurrent Neural Networks, arXiv preprint, 2016

TULIP Tool for Unsupervised Learning and Intelligence Processing

Advanced Data Analytics (G. Li)

9

More than images

Output ●●●●●●●●●●

Hidden Layer ●●●●●●●●●●

Hidden Layer ●●●●●●●●●●

Hidden Layer ●●●●●●●●●●

Input ●●●●●●●●●●

- Audio: Aaron van den Oord, Sander Dieleman, Heiga Zen, Karen Simonyan, Oriol Vinyals, Alex Graves, Nal Kalchbrenner, Andrew Senior, Koray Kavukcuoglu, WaveNet: A Generative Model for Raw Audio, arXiv preprint, 2016
- Video: Nal Kalchbrenner, Aaron van den Oord, Karen Simonyan, Ivo Danihelka, Oriol Vinyals, Alex Graves, Koray Kavukcuoglu, Video Pixel Networks , arXiv preprint, 2016

TULIP Tool for Unsupervised Learning and Intelligence Processing

Advanced Data Analytics (G. Li)

10

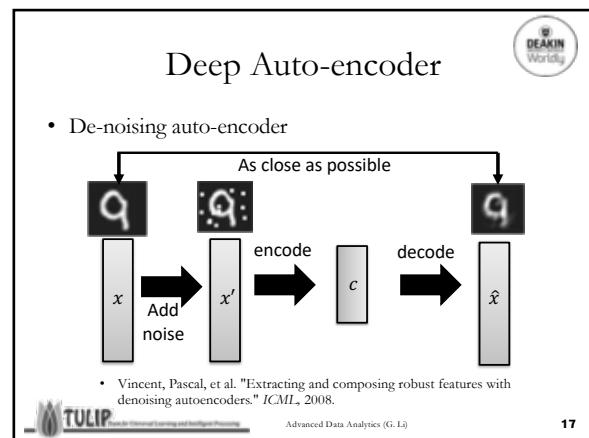
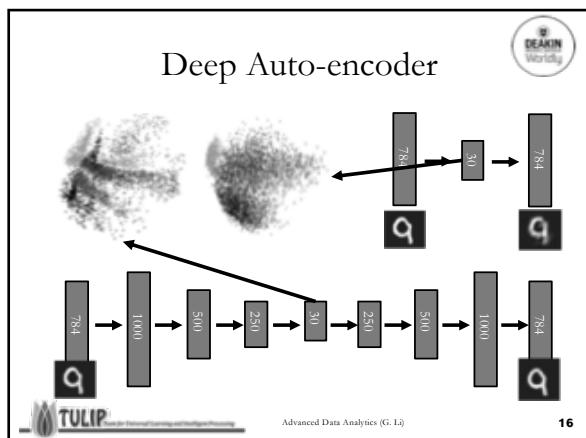
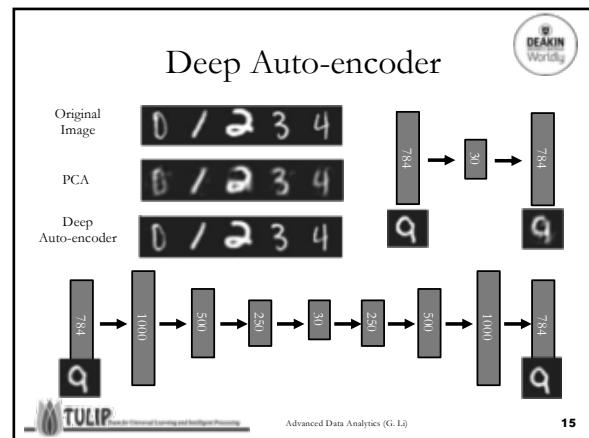
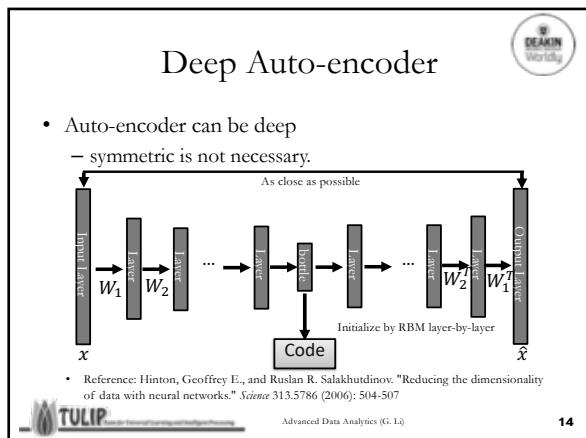
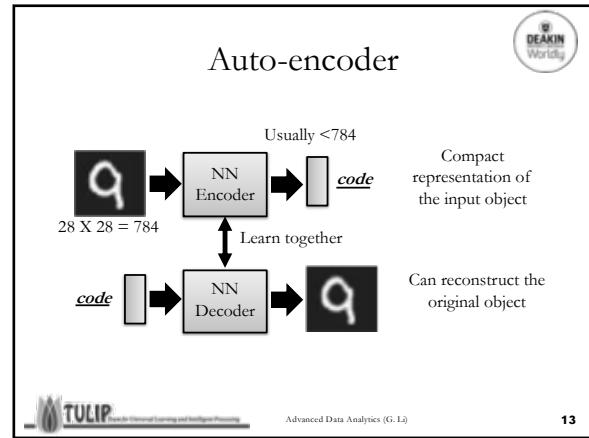
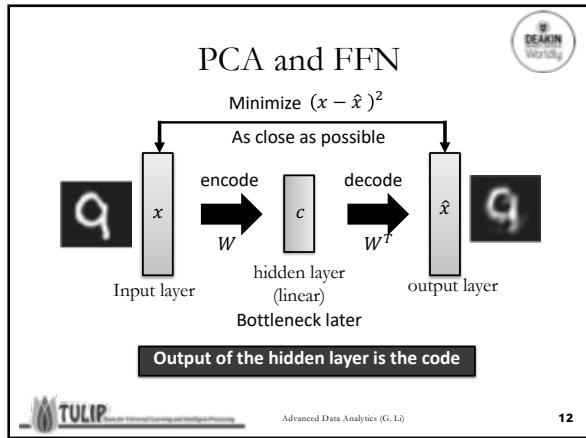
Auto-Encoder

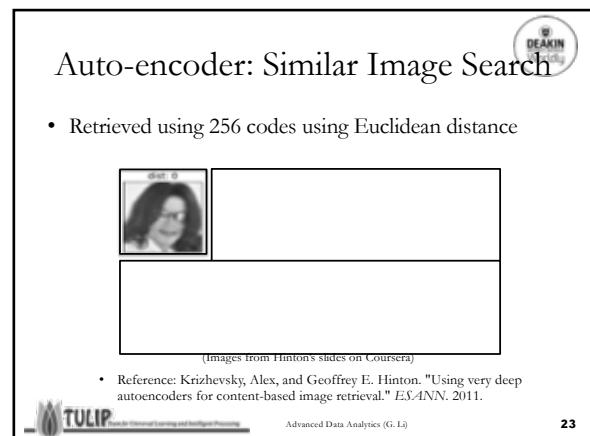
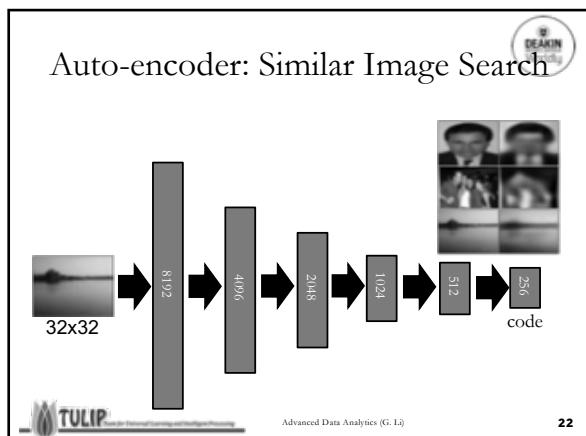
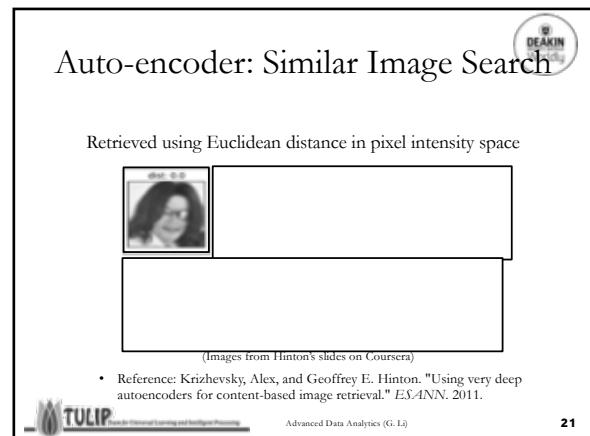
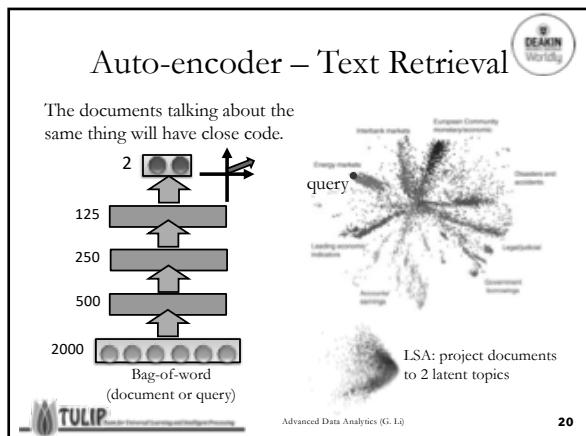
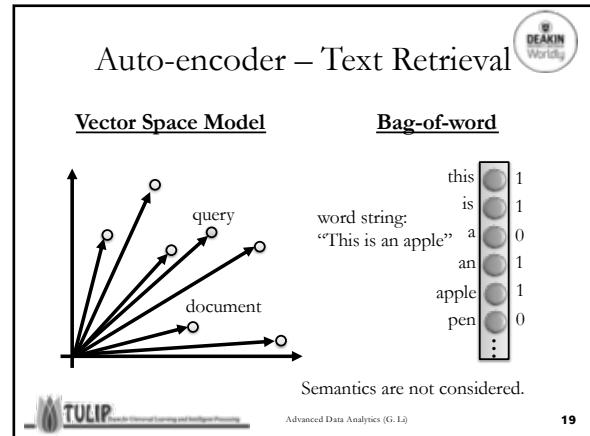
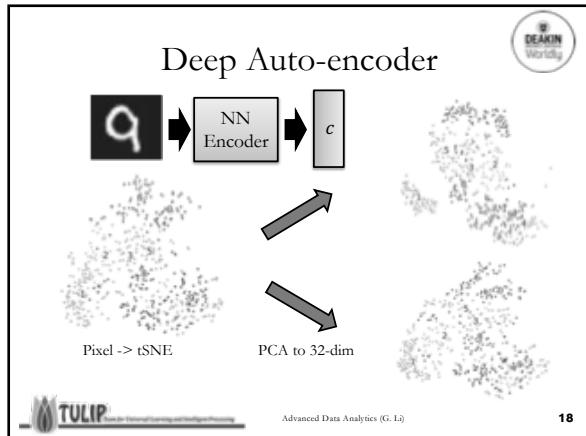
- Auto-Encoder
- Deep Auto-Encoder
- Auto-Encoder Applications

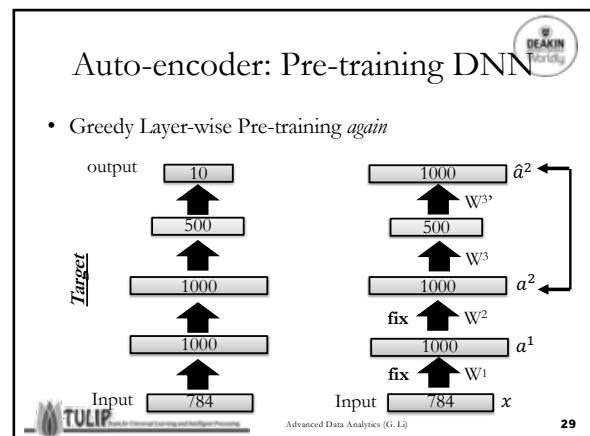
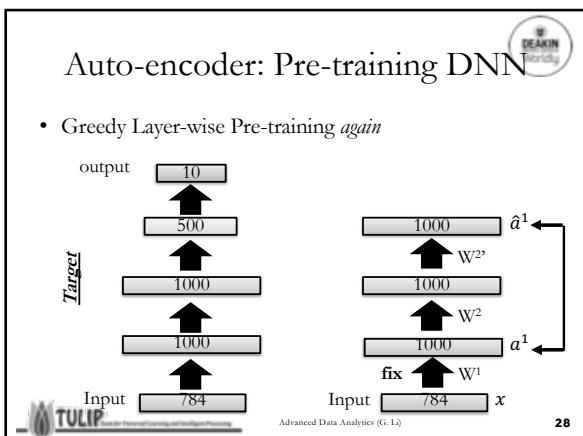
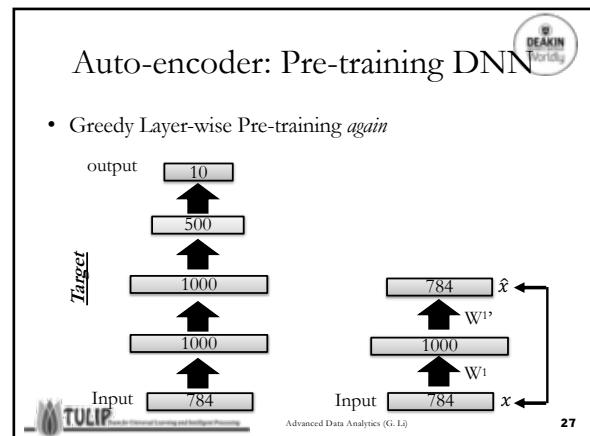
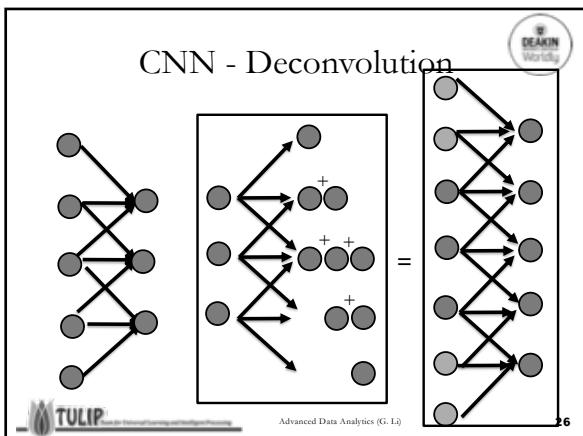
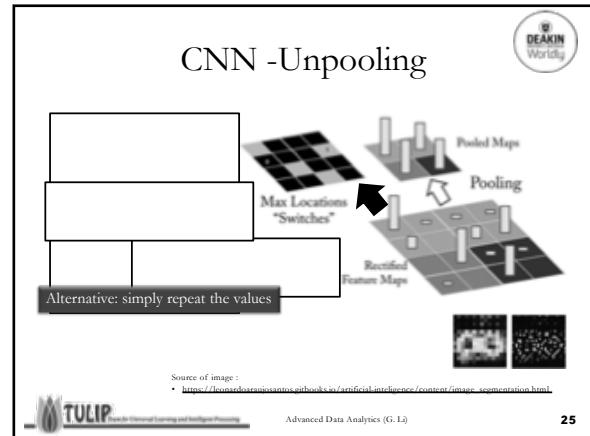
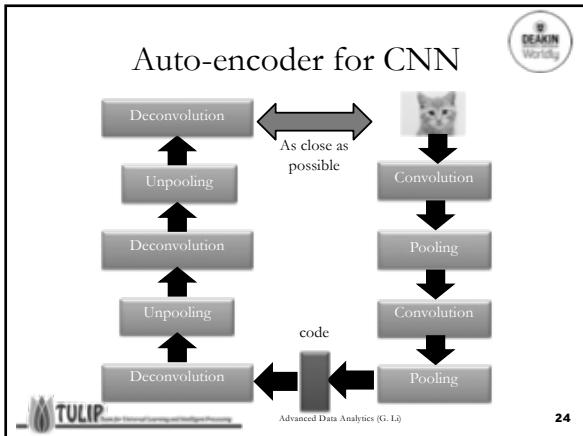
TULIP Tool for Unsupervised Learning and Intelligence Processing

Advanced Data Analytics (G. Li)

11

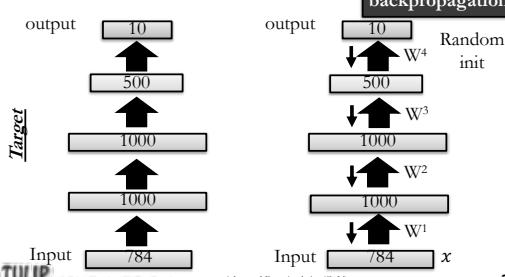






Auto-encoder: Pre-training DNN

- Greedy Layer-wise Pre-training *again*



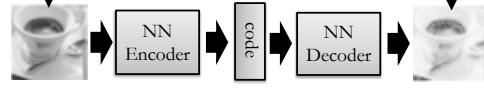
Advanced Data Analytics (G. Li)

30

Auto-encoder: Generating New Things

Find-tune by backpropagation

As close as possible



Randomly generate a vector as code

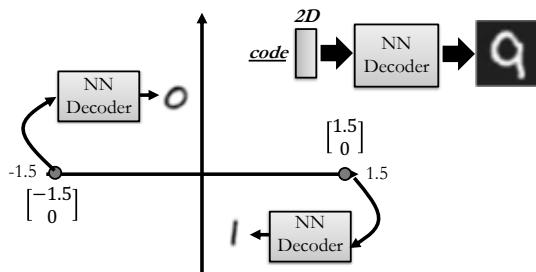
Image ?

TULIP

Advanced Data Analytics (G. Li)

31

Auto-encoder: Generating New Things

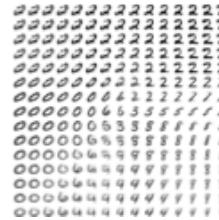
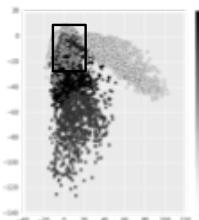


Advanced Data Analytics (G. Li)

32

Auto-encoder: Generating New Things

- Can we use decoder to generate something?



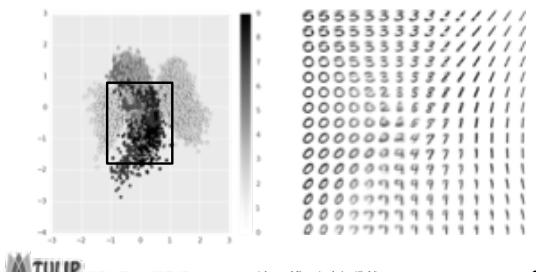
TULIP

Advanced Data Analytics (G. Li)

33

Auto-encoder: Generating New Things

- Can we use decoder to generate something?



Advanced Data Analytics (G. Li)

34

Variational Auto-Encoder

- Why VAE
- How VAE



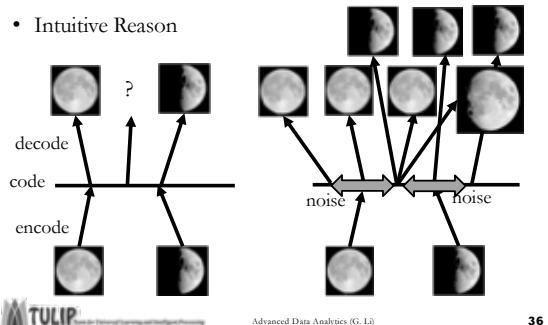
TULIP

Advanced Data Analytics (G. Li)

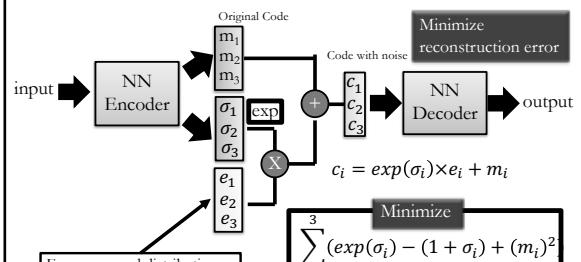
35

Why VAE?

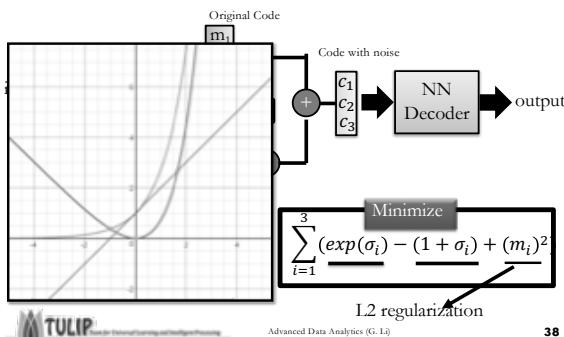
- Intuitive Reason



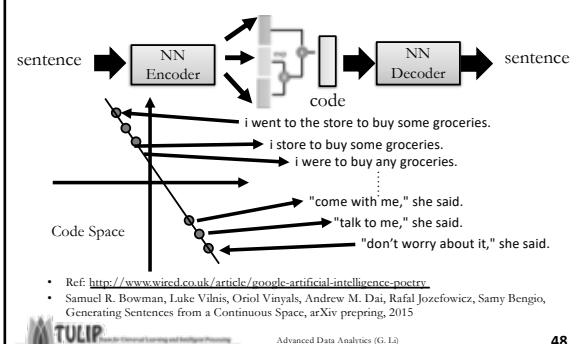
Variational Auto-Encoder



Variational Auto-Encoder

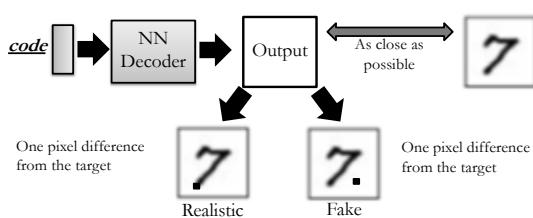


VAE Writing



Problems of VAE

- It does not really try to simulate real images



- VAE may just memorize the existing images, instead of generating new images

Generative Adversarial Network

- Why GAN
- How GAN



Cifar-10

- Which one is machine-generated?

Ref: <https://openai.com/blog/generative-models/>

TULIP: Data for Unsupervised Learning and Intelligence Processing

Advanced Data Analytics (G. Li)

51

Yann LeCun's comment

What are some recent and potentially upcoming breakthroughs in deep learning?

Yann LeCun, Director of AI Research at Facebook and Professor at NYU
Written Jul 29 - Upvoted by Joaquin Quiñonero Candela, Director Applied Machine Learning at Facebook and Huang Xiao

Adversarial training is the coolest thing since sliced bread.

I've listed a bunch of relevant papers in a previous answer.

Expect more impressive results with this technique in the coming years.

What's missing at the moment is a good understanding of it so we can make it work reliably. It's very finicky. Sort of like ConvNet were in the 1990s, when I had the reputation of being the only person who could make them work (which wasn't true).

Ref: <https://www.quora.com/What-are-some-recent-and-potentially-upcoming-breakthroughs-in-deep-learning>

TULIP: Data for Unsupervised Learning and Intelligence Processing

Advanced Data Analytics (G. Li)

52

Yann LeCun's comment

What are some recent and potentially upcoming breakthroughs in deep learning?

Yann LeCun, Director of AI Research at Facebook and Professor at NYU
Written Jul 29 - Upvoted by Joaquin Quiñonero Candela, Director Applied Machine Learning at Facebook and Nikhil Garg. I lead a team of Quora engineers working on ML/NLP problems

The most important one, in my opinion, is adversarial training (also called GAN for Generative Adversarial Networks). This is an idea that was originally proposed by Ian Goodfellow when he was a student with Yoshua Bengio at the University of Montreal (he since moved to Google Brain and recently to OpenAI).

This, and the variations that are now being proposed is the most interesting idea in the last 10 years in ML, in my opinion.

Ref: <https://www.quora.com/What-are-some-recent-and-potentially-upcoming-breakthroughs-in-deep-learning>

TULIP: Data for Unsupervised Learning and Intelligence Processing

Advanced Data Analytics (G. Li)

53

Evolution

Kallima inachus

Brown

veins

.....

Butterflies are not brown

Butterflies do not have veins

.....

TULIP: Data for Unsupervised Learning and Intelligence Processing

Advanced Data Analytics (G. Li)

54

The evolution of generation

NN Generator v1 → NN Generator v2 → NN Generator v3

Discriminator v1 → Discriminator v2 → Discriminator v3

Real images:

TULIP: Data for Unsupervised Learning and Intelligence Processing

Advanced Data Analytics (G. Li)

55

Basic Idea of GAN

- The data we want to generate has a distribution $P_{data}(x)$

$P_{data}(x)$

High Probability

Image Space

Low Probability

TULIP: Data for Unsupervised Learning and Intelligence Processing

Advanced Data Analytics (G. Li)

56

Basic Idea of GAN

- A generator G is a network. The network defines a probability distribution.

Normal Distribution z → generator G → $P_G(x)$
 $Q_x = G(z)$ → $P_{data}(x)$
 It is difficult to compute $P_G(x)$
 As close as possible
 We do not know what the distribution looks like.

Advanced <https://blog.openai.com/generative-models> 57

Basic Idea of GAN

Normal Distribution → NN Generator v1 → $P_G(x)$
 $P_G(x)$ vs $P_{data}(x)$
 image → Discriminator v1 → 1/0
 It can be proved that the loss the discriminator related to JS divergence

Advanced Data Analytics (G. Li) 58

Basic Idea of GAN

- Next step**
 - Updating the parameters of generator
 - To minimize the JS divergence

The output be classified as “real” (as close to 1 as possible)

Generator + Discriminator = a network
 Using gradient descent to update the parameters in the generator, but fix the discriminator

Normal Distribution → NN Generator v2 → Generator output → Discriminator v1 → 1.0 / 0.13

Advanced Data Analytics (G. Li) 59

GAN- DCGAN

Source of images: <https://zhuanlan.zhihu.com/p/24767059>
 DCGAN: <https://github.com/carpedm20/DCGAN-tensorflow>

Advanced Data Analytics (G. Li) 60

GAN- DCGAN

100 rounds

Advanced Data Analytics (G. Li) 61

GAN- DCGAN

1000 rounds

Advanced Data Analytics (G. Li) 62

GAN- DCGAN

2000 rounds



63

TULIP Tool for Unsupervised Learning and Intelligence Processing

Advanced Data Analytics (G. Li)

GAN- DCGAN

5000 rounds



64

TULIP Tool for Unsupervised Learning and Intelligence Processing

Advanced Data Analytics (G. Li)

GAN- DCGAN

10,000 rounds



65

TULIP Tool for Unsupervised Learning and Intelligence Processing

Advanced Data Analytics (G. Li)

GAN- DCGAN

20,000 rounds



66

TULIP Tool for Unsupervised Learning and Intelligence Processing

Advanced Data Analytics (G. Li)

GAN- DCGAN

50,000 rounds



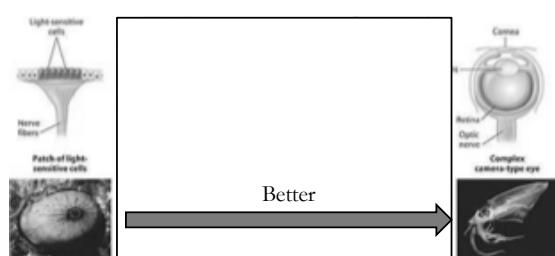
67

TULIP Tool for Unsupervised Learning and Intelligence Processing

Advanced Data Analytics (G. Li)

<http://www.guokr.com/post/772890/>

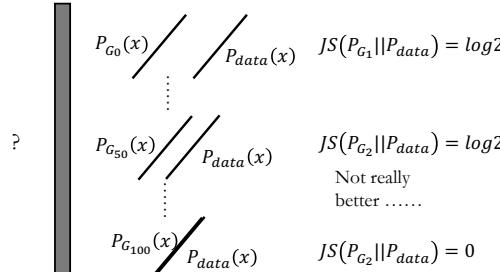
Why GAN is hard to train?



Advanced Data Analytics (G. Li)

68

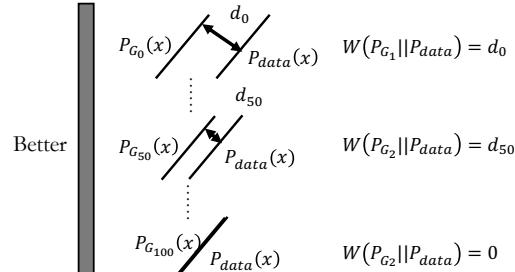
Why GAN is hard to train?



Advanced Data Analytics (G. Li)

69

WGAN (Wasserstein distance)



Advanced Data Analytics (G. Li)

70

Moving on the code space

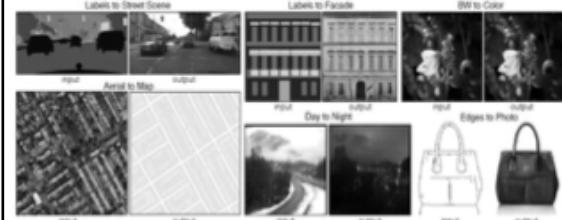


Alec Radford, Luke Metz, Soumith Chintala, Unsupervised Representation Learning with Deep Convolutional Generative Adversarial Networks, ICLR, 2016

Advanced Data Analytics (G. Li)

71

Image-to-image Translation

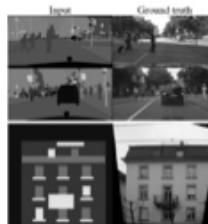


Philip Isola, Jun-Yan Zhu, Tinghui Zhou, Alexei A. Efros, "Image-to-Image Translation with Conditional Adversarial Networks," arXiv preprint, 2016

Advanced Data Analytics (G. Li)

72

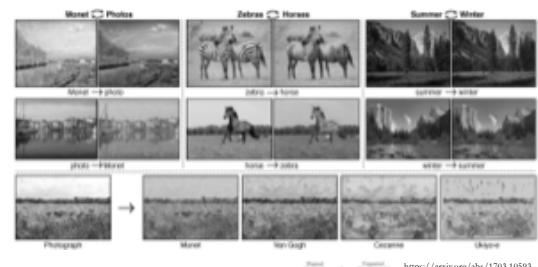
Image-to-image Translation - Results



Advanced Data Analytics (G. Li)

73

Cycle GAN



Advanced Data Analytics (G. Li)

74

Variants of GAN



Modifying the Optimization of GAN

fGAN
WGAN
Least-square GAN
Loss Sensitive GAN
Energy-based GAN
Boundary-seeking GAN
Unroll GAN
.....

Different Structure from the Original GAN

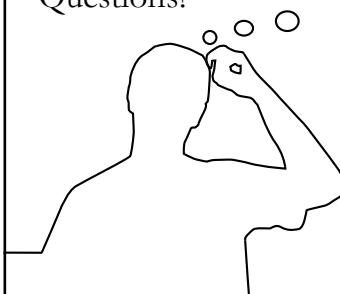
Conditional GAN
Semi-supervised GAN
InfoGAN
BiGAN
Cycle GAN
Disco GAN
VAE-GAN
.....



Advanced Data Analytics (G. Li)

75

Questions?



Advanced Data Analytics (G. Li)

76