

Spark for data science

→ Data is getting big & we need ways to handle it.

Lesson overview

- ① Big data ecosystem
- ② Data wrangling
- ③ Debugging & optimisation
- ④ Spark ML

The power of spark

Big data → lots of tasks (ecosystem)

↖ ? what is this

- Data that cannot be processed on a single machine.

→ important numbers

CPU .4 us ← "Brain" of pc
Memory 100 us ← "Scratch Paper"
Storage 16 ms ← "File cabinet"
Network 150 ms ← "outside world"

CPU IS FAST

problem: loading data takes longer than
time taken to process

↳ we'll also lose data when pc is
turned off.

↳ is also expensive

Storage is slower still, but it is cheaper

Network even slower

← tends to be bottleneck
when working w/ big data

Sparks ↓ data shuffling

When should you use spark?

↳ too big to process on 1 pc

∴ split the data & distribute it

↳ email?
 ↳ diff versions
 ↳ slow
 ↳ not controlled

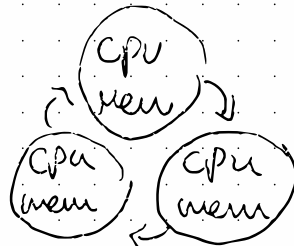
→ Do it to prevent "thrashing"

↳ loading data in
and out of memory

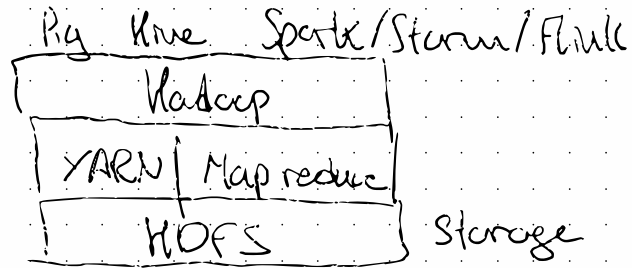
→ even if data is bigger than your
local memory, you can still load it
using spark to handle the data in
batches

↳ Distributed
computation

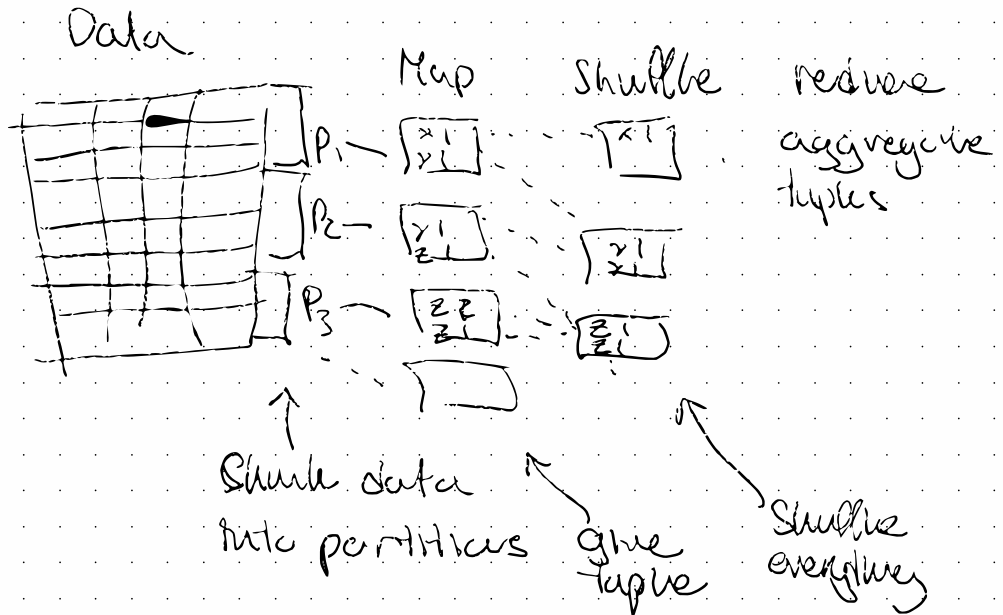
Parallel comp.
 ↔ CPU
 ↔ CPU
 ↔ CPU



Hadoop framework



Mapreduce



Data wrangling w/ spark

Spark ← functional programming
paradigm (Scala)

↖ access through
Python via API
(PySpark)

Why use functional programming?

(1) gets name from "functions" in
mathematics

Not the same as procedural prog.

ex. if you run a counter on an
iterator many times you will get
a different output.

(2) Functions should not access
global variables

(3) No changes made to input data
either.

(4) Chaining of functions together

How to make a spark script

- ↳ Parallelise dataset
- ↳ use a lot of map functions (w/ lambda func.)
- ↳ only runs at the last second

basic workflow:

```
import pyspark as ps
```

```
sc = ps.SparkContext(name)
```

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

```
rdd.map(func).collect()
```

↑
can be
a func or
lambda

↑ only runs
at this point

- ↳ select from rdd (resilient dist. dataset)
using common keywords

```
rdd.where  
rdd.select
```

? order doesn't matter

↳ or use `rdd.sql()`
for sql

Data formats

() CSV, JSON, XML, HTML, ...

Distributed Data Stores

() stores data across a cluster (HDFS)

· Adds redundancy so if 1 node dies then no data is lost

Spark Sessions

Spark
↓ Program

Spark Context

↳ Spark Conf in order to configure settings.

Imperative vs Declarative Programming



Python

"How?" ← Path to results



SQL

"What" ← Result

this means that when you tell SQL what you want to see, it'll know exactly how best to make it.

- () adv. does one thing well
- () disadv. only 1 thing done well.

Aggregations

done w/ `rdd.agg()`

↑ dict w/
`{ col : function }`?

rename to useable w/

`rdd.withColumnRenamed('agg(col)', <new>)`

Machine Learning w/ Spark

Libraries → Spark ML
→ Spark ML Lib ← slowly getting deprecated

Feature extraction

↳ information needs to be extracted from data to make features

c) Also needs to be normalised & put into the correct format.

Spark is a little weird...

↳ each row needs to become a vector to be used in functions

test f1 f2 f3
 <64t>

x	y	z
---	---	---

 → [x, y, z]
 ↘ ints of floats ↗ using vector assembler

General Workflow for ml

↳ wrangle data



turn to vec & label col



place into ml algorithm

train/test
← splits?

k-fold
cv?

Pipelines

Spark ml has similar pipeline capabilities to sklearn

- transformers ← feature extraction
- regressors ← ML algorithm

↳ pipeline (stages = <list of transformers>)
w/ est at the end

Model Selection & Tuning

param grid

↳ Has params to be searched

evaluator

↳ evaluates performance

↳ found in `pyspark.ml.tuning`