# Lecture 22:

## Data extraction

abstraction:

extract - combine

extract combine :

$$('k * 'k \longrightarrow order)$$

$$* ('a \longrightarrow ('k * 'v) \text{ Seq.seq})$$

$$* ('v * 'v \longrightarrow 'v)$$

$$* 'a \text{ Seq.seq}$$

$$\longrightarrow ('k, 'v) \text{ Dict.dict}$$

Comparison

eXtractor

Combiner

'a is a doc
(documents)

data

StringSequence

○  ○ ○ ○ ○ ○

this
is is
doc
1

○ ○ ○

this
is
doc
2

extractor

⟨ ("this", 1),
(is, 1),
(is, 1),
(doc, 1),
(1, 1) ⟩

extractor

⟨ (this, 1)
(is, 1)
(doc, 1),
(2, 1) ⟩

# Dictionary where:

$$\{ \text{"this"} \sim 2,$$
$$\text{"is"} \sim 3,$$
$$\text{"doc"} \sim 2,$$
$$\text{"I"} \sim 1,$$
$$\text{"2"} \sim 1 \}$$

key                    value

```sml
fun wordcount (docs: string Seq.seq): (string, int)
                                        Dict.dict =
  extract combine
    ( String.compare,
```

extractor
```sml
      fn doc: string => Seq.map ( fn w => (w, 1),
                                  words doc ))
```

combiner
```sml
      fn (x, y) => x+y,
```

```sml
      docs )
```

## MAP-reducable type

| lists | trees | seq |
| --- | --- | --- |
| map | map | map |
| reduce | reduce | reduce |

| files |
| --- |

map +
reduce ?

networked
cluster
of
computers

map +
reduce

P

Cache

RAM

hard
drive

files

sequences
trees
lists

born that sim

Six computer nodes connected in a distributed network. Each node contains: P, Cache, RAM, and a hard drive.

Node 1 hard drive: file 1
Node 2 hard drive: file 2
Node 3 hard drive: file 3, file 4
Node 4 hard drive: file 5
Node 5 hard drive: file 6

```
signature MAP_REDUCE =
sig
    type 'a mapreducable

    val map: ('a -> 'b) *
             ('a mapreducable)
        -> 'b mapreducable

    val reduce: ('a * 'a -> 'a)
             * 'a
             * 'a mapreducable
        -> 'a
end
```

short summary of lots of data

reduce ( +,

0

map (fn x => x+1, <1,4,8...>

|=> reduce ( +, 0, <2,5,9 ----- =>)

|=> 2+5+9

```
Signature MAP_REDUCE =
Sig
    type 'a mapreducable

    val mapreduce : ('a -> 'b)
                    * 'b
                    * ('b * 'b -> 'b)
                    * ('a mapreducable)
                    -> 'b

end
```

$$\text{mapreduce} \left( \text{fn } x \Rightarrow x + 1, \right.$$
$$0,$$
$$+,$$
$$\left. \langle 1, 4, 8, - - - - - - \rangle \right)$$

$$=$$
$$2 + 5 + 9 + - - - - -$$

```sml
structure SeqMR :> MAP_REDUCE =
struct

    type 'a mapreducable =
        'a Seq.seq

    fun mapreduce (f, e, n, s) =
        Seq.reduce(n, e,
                   Seq.map(f, s))
end
```

Structure File MR: MAP_REDUCE

data is stored
in a file

+ mapreduce over it!

training
+ test