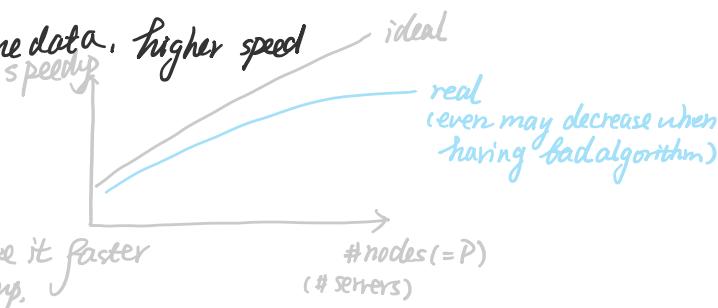
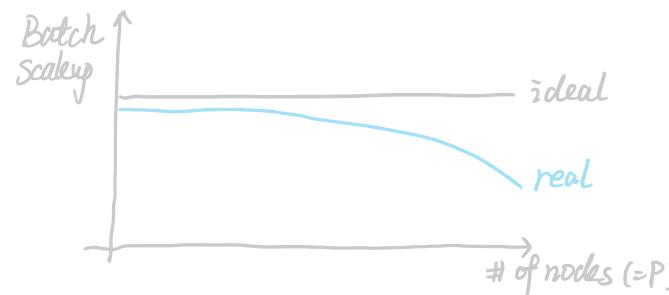


# Spark

Wednesday, May 15, 2019 1:30 PM

- Why parallelism:
  - Multi-cores
  - big data too large to fit in main memory
    - distributed query processing on 100x-1000x servers
- Performance Matrix for DBMS
  - Nodes: processors, computers

- speed up: more nodes, same data, higher speed
  - Linear vs Non-linear
  - there's no way to make it faster than linear speedup.
- scale up: more nodes, same speed, more data
  - Linear vs Non-linear scaleup

- Reason for sub-linear:

{ start up cost cost of an starting operation on many nodes  
 interference contention for resources between nodes  
 skew slowest node becomes the bottleneck

- **Spark** A case study of the MapReduce Programming paradigm.

- Distributed processing over HDFS.
- Supports interfaces in Java, Scala, Python.
- Difference from MapReduce
  - { Multiple steps including iterations
  - { Store intermediate result in main memory
  - Closer to relational algebra
- A spark program
  - { Transformation (map, join, reduce) **Lazy**
  - { Actions (count, reduce, save) **Eager**
    - Lazy: operator tree is constructed in memory
    - Eager: operators are executed immediately
- Collections:
  - **RDD<T>** = an RDD collection of type T
    - **not nested, distributed.**
    - **executed in parallel**
    - **Recoverable** via **lineage**

- $\text{Seq} < T >$  = a sequence

- local to one server, may be nested
- done sequentially

e.g.,  $s = \text{SparkSession.builder()...getOrCreate();}$

```
sqlErrors = s.read().textFile("hdfs://logfile.log")
    .filter(l -> l.startsWith("ERROR"))
    .filter(l -> l.contains("sqlite"))
    .collect();
```

action

↑ **Call chaining style.**

write the functions as a chain

- Anonymous functions / **Lambda Expressions**

```
errors = lines.filter(l -> l.startsWith("Error"));
```



```
class FilterFn implements Function<Row, Boolean>{
```

```
    Boolean call(Row r)
```

```
    { return l.startsWith("Error"); }
```

}

```
errors = lines.filter(new FilterFn());
```

- **Fault Tolerance**

When a job is executed on many servers, the probability of failure is high.

$$P(\text{fail}) = \sum P_i(\text{fail})$$

is the sum of  $P$  of all servers.

◦ soln:

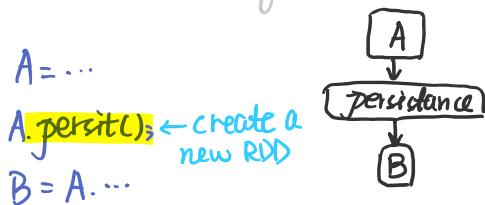
- **Parallel database system** restart the whole thing (Expensive)

- **MapReduce** write everything to disk, redo. (Slow)

- **Spark** redo only what's needed. (efficient)

- **RDD = Resilient Distributed Dataset**
  - *Distributed, immutable, and records its lineage.*
  - *Lineage = expression that says how that relation was computed = a relational algebra plan*
- **Spark stores intermediate results as RDD.**
- *If a server crashes, its RDD in main memory is lost. However, the driver (= master node) knows the lineage, and will simply recompute the lost partition of the RDD.*

- **Persistence** : a kind of immediate materialization



Transformations:	
<code>map(f : T -&gt; U):</code>	<code>RDD&lt;T&gt; -&gt; RDD&lt;U&gt;</code> <i>perform f on each T.</i>
<code>flatMap(f: T -&gt; Seq(U)):</code>	<code>RDD&lt;T&gt; -&gt; RDD&lt;U&gt;</code> <i>perform f on T, get a sequence, then flat it</i>
<code>filter(f:T-&gt;Bool):</code>	<code>RDD&lt;T&gt; -&gt; RDD&lt;T&gt;</code>
<code>groupByKey():</code>	<code>RDD&lt;(K,V)&gt; -&gt; RDD&lt;(K,Seq[V])&gt;</code>
<code>reduceByKey(F:(V,V)-&gt; V):</code>	<code>RDD&lt;(K,V)&gt; -&gt; RDD&lt;(K,V)&gt;</code>
<code>union():</code>	<code>(RDD&lt;T&gt;,RDD&lt;T&gt;) -&gt; RDD&lt;T&gt;</code>
<code>join():</code>	<code>(RDD&lt;(K,V)&gt;,RDD&lt;(K,W)&gt;) -&gt; RDD&lt;(K,(V,W))&gt;</code>
<code>cogroup():</code>	<code>(RDD&lt;(K,V)&gt;,RDD&lt;(K,W)&gt;)-&gt; RDD&lt;(K,(Seq&lt;V&gt;,Seq&lt;W&gt;))&gt;</code>
<code>crossProduct():</code>	<code>(RDD&lt;T&gt;,RDD&lt;U&gt;) -&gt; RDD&lt;(T,U)&gt;</code>

Actions:	
<code>count():</code>	<code>RDD&lt;T&gt; -&gt; Long</code>
<code>collect():</code>	<code>RDD&lt;T&gt; -&gt; Seq&lt;T&gt;</code>

<code>reduce(t:(I,I)-&gt;I):</code>	RDD<I> -> I
<code>save(path:String):</code>	Outputs RDD to a storage system e.g., HDFS

- Three Java-Spark APIs

- RDD `javaRDD<T>`

- basically untyped
- distributed, main memory

- Data frames `DataSet<Row>`

- Row: a record, **dynamically typed.**
- Distributed, main memory, or external (e.g. SQL)

- Datasets `Dataset<Person>`

- `<Person>`: user defined type
- distributed, main memory (not external)

- **Data frames**

- like RDD, also **immutable** distributed collection of data.
- Organized into **named columns** rather than individual objects  
(like a relation)
- Similar API as RDDs with additional methods.

- **Dataset**

- Similar to DataFrames, except that the element must be typed.
- Can detect errors during compilation time

## Datasets API: Sample Methods

- Functional API
  - `agg(Column expr, Column... exprs)`  
Aggregates on the entire Dataset without groups.
  - `groupBy(String col1, String... cols)`  
Groups the Dataset using the specified columns, so that we can run aggregation on them.
  - `join(Dataset<?> right)`  
Join with another DataFrame.
  - `orderBy(Column... sortExprs)`  
Returns a new Dataset sorted by the given expressions.
  - `select(Column... cols)`  
Selects a set of column based expressions.
- “SQL” API
  - `SparkSession.sql("select * from R");`