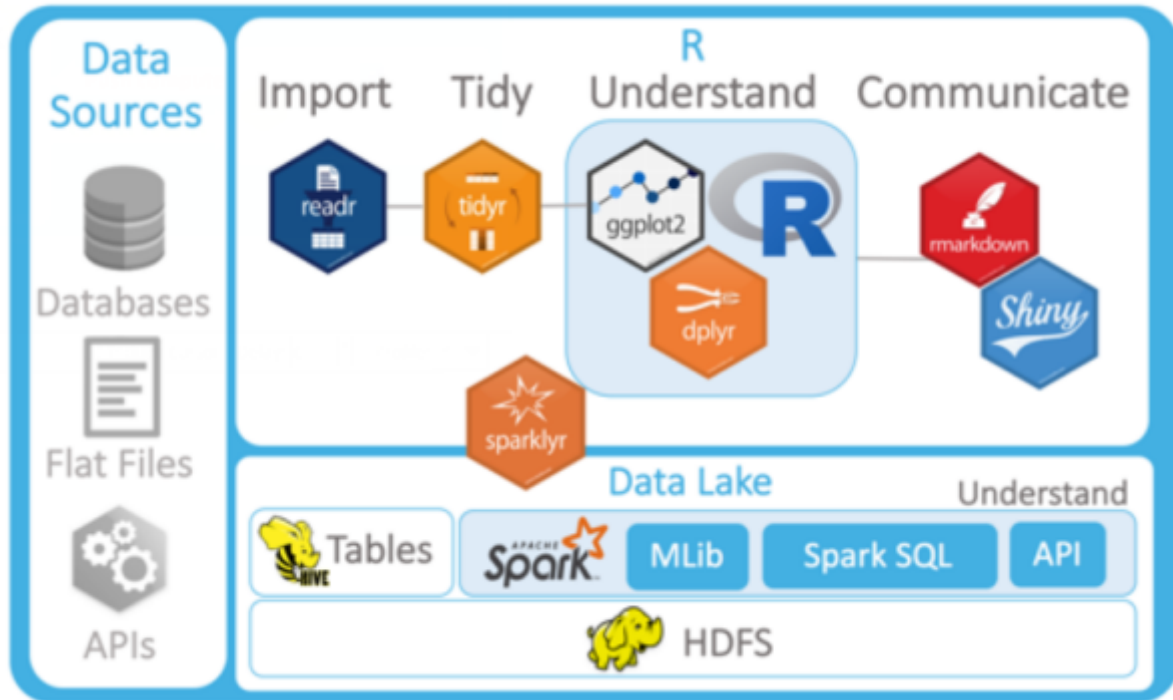


# Big data with Modern R & Spark

"SEMINAR: *Análisis de Big data con Tidyverse y Spark: uso en estadística pública*"

By Xavier de Pedro Puente, Ph.D.

Senior Technician at the Barcelona City Council.



Wednesday, June 12, 2024. 16:00-17:00h + 30' questions

Within the context of the postgraduate course on

**"Data Science. Applications to Biology and Medicine with Python and R"**

at University of Barcelona (IL3). 2024.

## Outline

1. About me
2. Barcelona City Council Case
3. Modern R: Tidyverse
4. Modern R & Big Data
5. Spark (Apache) & Sparklyr (RStudio)
6. Sparklyr - next steps
7. Alternative approach
8. Former hands-on exercise nowadays with sparklyr

# (1) About me

Xavier de Pedro Puente, Ph.D. [xavier.depedro \(a\) seeds4c.org](mailto:xavier.depedro@seeds4c.org)

- **Academics:**

- Degree in Biology (UB<sup>[1]</sup>)
- Ph.D. in Ecology (UB<sup>[2]</sup>)
- Postgraduate in Bioinformatics (UOC<sup>[3]</sup>)

- **Current Work:**

- Senior technician at **Climate Change and Sustainability Office** (Barcelona City Council<sup>[4]</sup>)

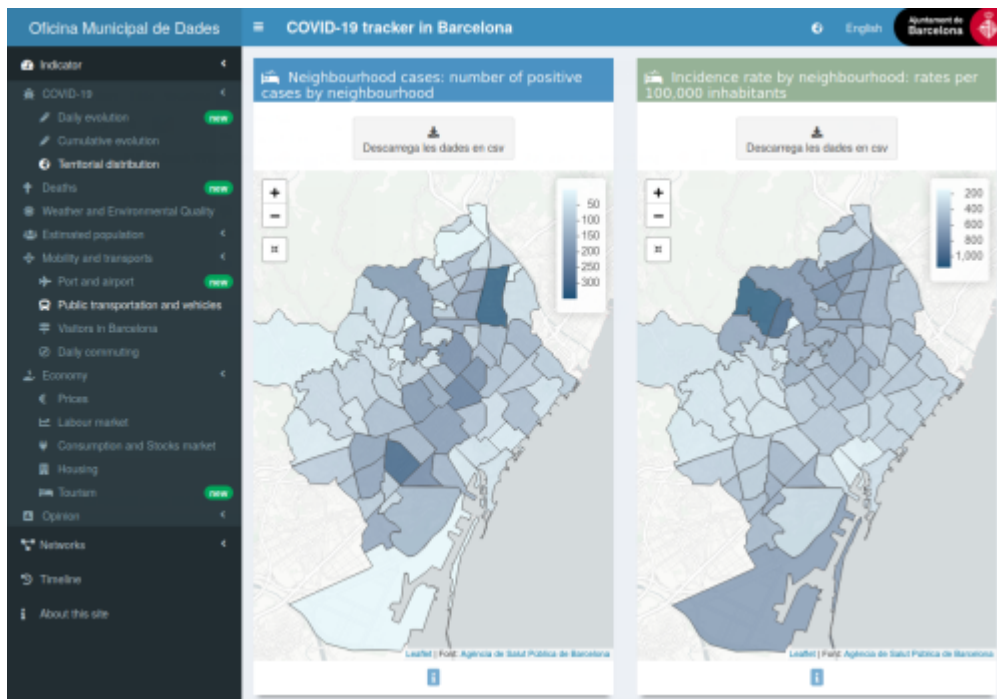
- **Past (related) Work:**

- Senior technician at Development Division, (IMI<sup>[5]</sup>, Barcelona City Council<sup>[6]</sup>)
- Senior technician at Municipal Data Office (OMD<sup>[7]</sup>, Barcelona City Council<sup>[8]</sup>)
- Bioinformatics technician (UEB<sup>[9]</sup>, VHIR<sup>[10]</sup>)
- Systems administrator (UEB<sup>[11]</sup>, VHIR<sup>[12]</sup>)

**Disclaimer** The views expressed in this presentation are those of the author and do not necessarily represent the views and policies of the Barcelona City Council. Any mention of trade names, products or services does not imply an endorsement by the Barcelona City Council unless explicitly stated as such.

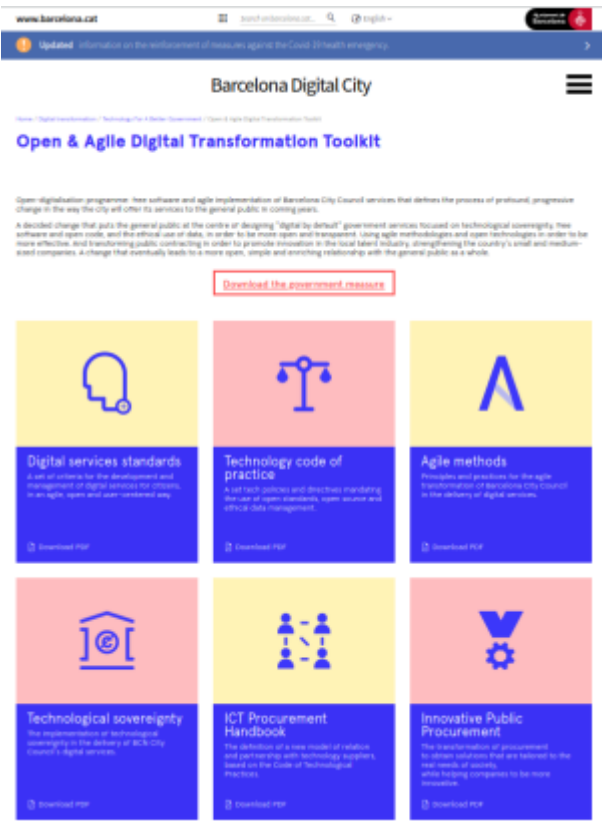
## (2) Barcelona City Council Case

- 2018: OMD - Municipal Data Office opened<sup>[13]</sup> (details<sup>[14]</sup>):  
Management, quality, governance and use of data controlled and/or stored by Barcelona City Council and all of its associated bodies (both public and private).
- Public Statistics<sup>[15]</sup> portal
- Recent product<sup>[16]</sup>: Covid-19 in BCN monitoring with an R Shiny App<sup>[17]</sup>



## Barcelona City Council: Public money - public code

- Barcelona: first city to join the Free Software Foundation, Public Money, Public Code<sup>[18]</sup> campaign
- One of the case studies<sup>[19]</sup> of the use of open-source software and open code to democratise cities.



There is a Government policy for open source and agile methodologies<sup>[20]</sup>

## Barcelona City Council: CityOS

City Council's infrastructure based on open-source Big Data technology

Video<sup>[22]</sup>

**City OS:** internal data management, known as "Data Lake"<sup>[21]</sup>

# OPEN SOURCE URBAN PLATFORM

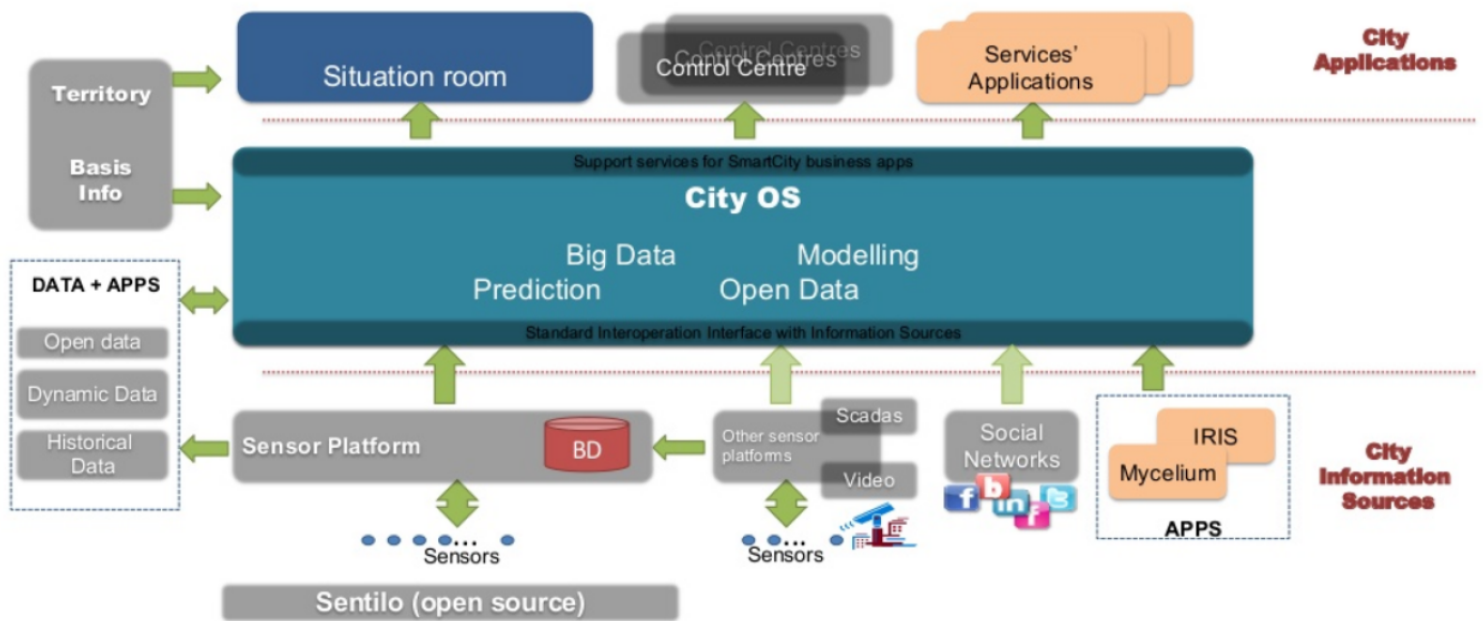
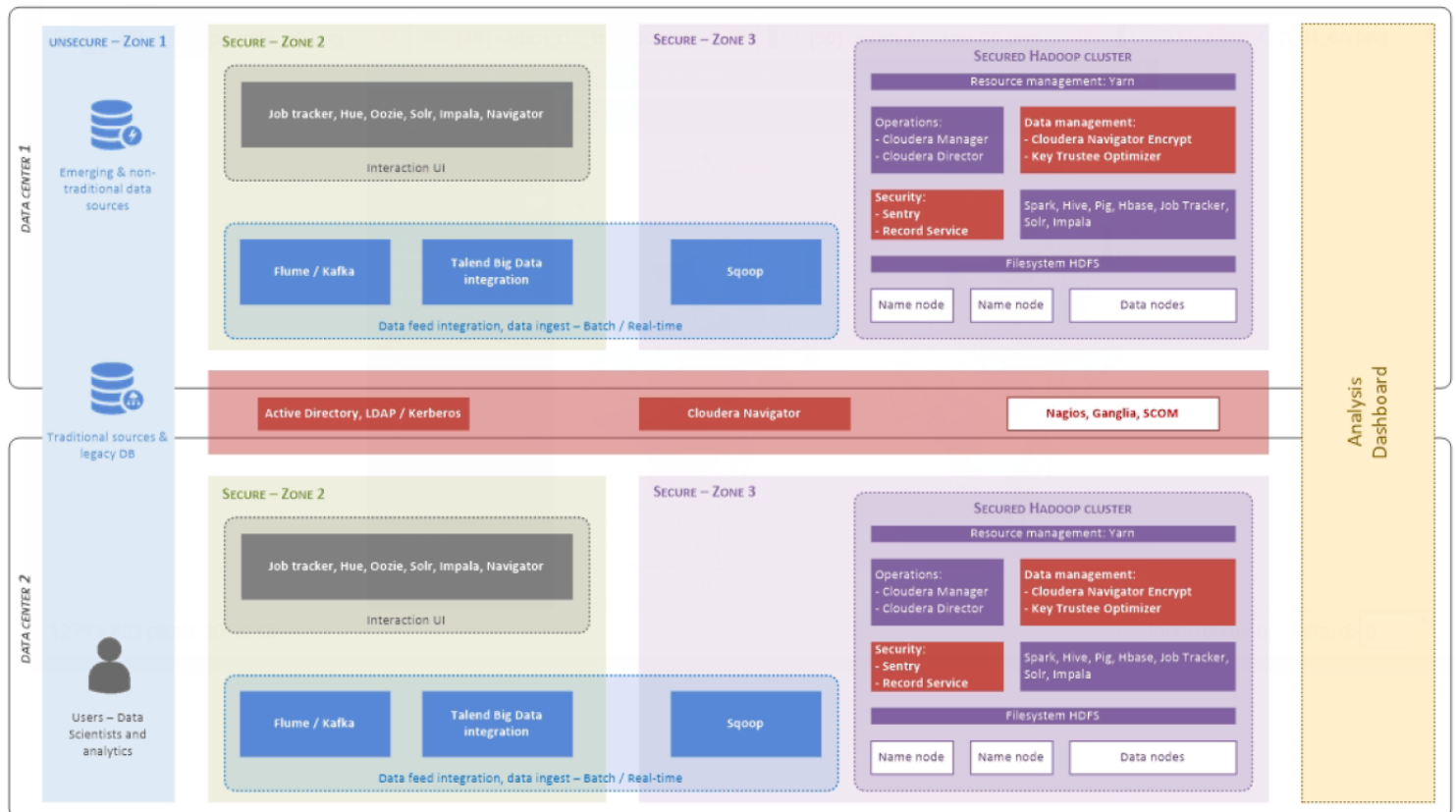


Image from Francesca Bria<sup>[23]</sup> (Barcelona Digital City Roadmap 2017-2020)

## Barcelona City Council: CityOS Technologies

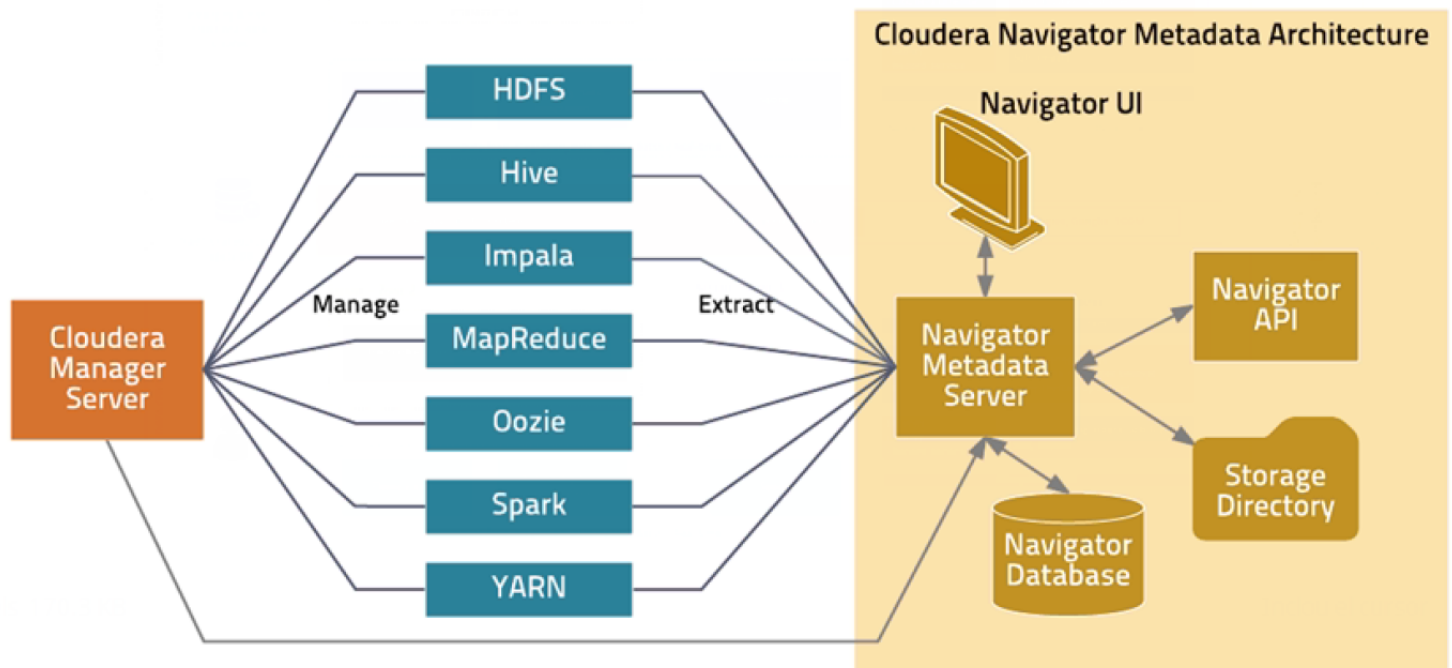
CityOS technologies: GNU/Linux, CentOS, Cloudera, Activiti, Talend, Protégé, **R**, Zabbix, Nagios, Ganglia



Source: [https://github.com/AjuntamentdeBarcelona/CityOS\\_AjBCN](https://github.com/AjuntamentdeBarcelona/CityOS_AjBCN)<sup>[24]</sup> - Image from J. Berdonces a Github<sup>[25]</sup>

# Barcelona City Council: CityOS - Cloudera Manager

Cloudera with Hadoop File System, Hue, HBase, Hive, Impala, Oozie, **Spark**, Yarn, Kafka, ...



Source: [https://github.com/AjuntamentdeBarcelona/CityOS\\_AjBCN](https://github.com/AjuntamentdeBarcelona/CityOS_AjBCN)<sup>[26]</sup> - Image from J. Berdonces a Github<sup>[27]</sup>

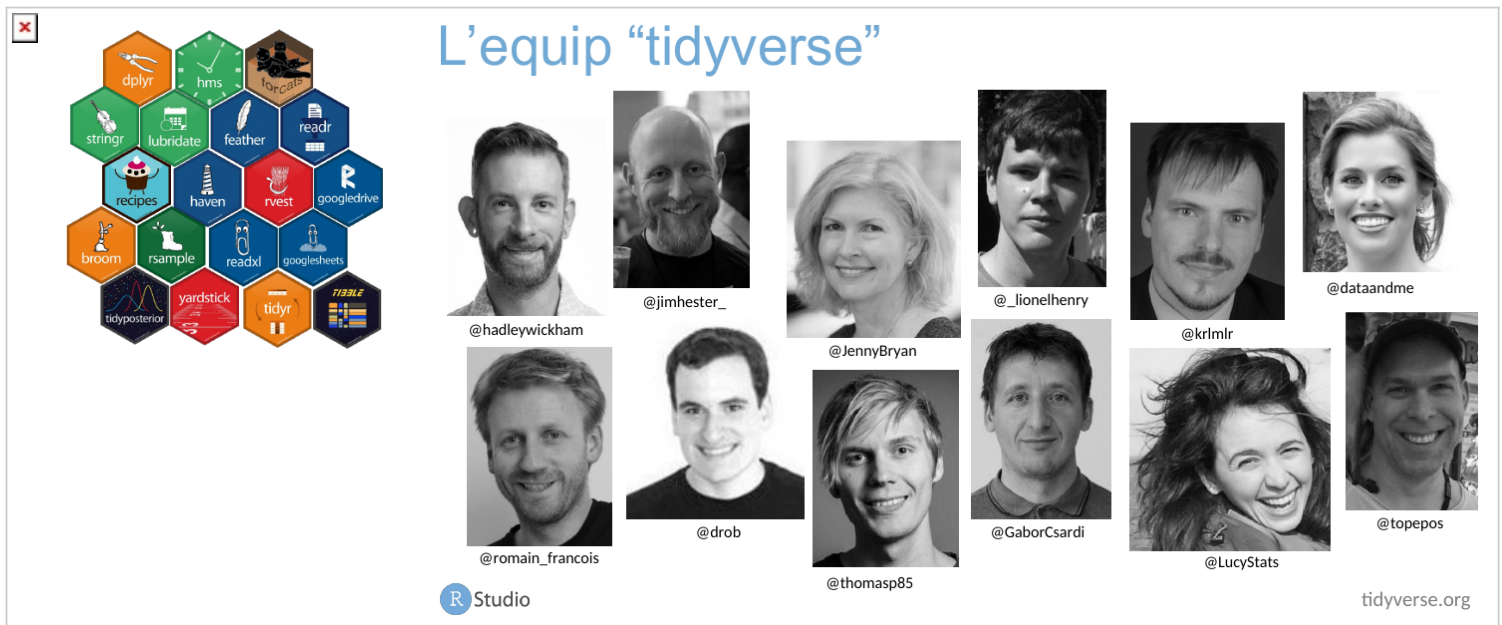
## (3) From Base R to "Modern R": Tidyverse

- 2000: First stable beta version (v1.0) released.
- ... "Base R" ...
- 2016: Tidyverse<sup>[28]</sup> 1.0.0 package released on CRAN<sup>[29]</sup>
- 2016: "Introduction to Modern R - r4stats.com"<sup>[30]</sup>
- 2017: "Martin Hadley on R and the Modern R Ecosystem - InfoQ Podcasts"<sup>[31]</sup>
- 2017: "Modern Data Science with R - Chapman & Hall/CRC"<sup>[32]</sup>
- 2019: "Modern R with the tidyverse - Econometrics and Free Software"<sup>[33]</sup> (blog)
- 2019: "Modern R and the Tidyverse - Data Science Workshops"<sup>[34]</sup>
- 2020: "Modern R: Welcome to the tidyverse"<sup>[35]</sup> (video tutorial)
- 2020: Statistical Inference via Data Science: ModernDive into R & Tidyverse<sup>[36]</sup>
- 2022: Modern R with the Tidyverse<sup>[37]</sup>



## Modern R (Tidyverse)

Workflow, Packages, People/Community



Derived from here<sup>[38]</sup> & here<sup>[39]</sup>

## Modern R (Tidyverse) principles

### 1. Main structures are ordered data

- Each variable is saved in its own column.
- Each observation is saved in its own row.
- Each "type" of observation stored in a single table

```
## country 2011 2012 2013
## 1 FR 7000 6900 7000
## 2 DE 5800 6000 6200
## 3 US 15000 14000 13000
```

```
cases %>%
  pivot_longer(
    cols=!country,
    names_to="year",
    values_to="n"
  )
```

```
## country year n
## 1 FR 2011 7000
## 2 DE 2011 5800
## 3 US 2011 15000
## 4 FR 2012 6900
## 5 DE 2012 6000
## 6 US 2012 14000
## 7 FR 2013 7000
## 8 DE 2013 6200
## 9 US 2013 13000
```

~~gather(cases, "year", "n", 2:4)~~

### 2. Each function represents one step

### 3. Functions are combined with the pipe

**operator** %>%

- or with the earlier pipe<sup>[40]</sup> from base R: |>

### 4. Each step is a query or command

Derived from here<sup>[41]</sup>, here<sup>[42]</sup> & here<sup>[43]</sup>

## (4) Modern R & Big Data

# Data > RAM = Problem



# ✈ Airlines Data Set

Arrival and departure details for all commercial flights in US between October 1987 and April 2008.

120,000,000 records. **12 GB**

[stat-computing.org/dataexpo/2009/](http://stat-computing.org/dataexpo/2009/)



Data does not  
fit in memory

From an RStudio seminar<sup>[44]</sup> by Garrett Grolemond<sup>[45]</sup>

## 3 classes of Big Data Problems

### Class 1. Extract Data

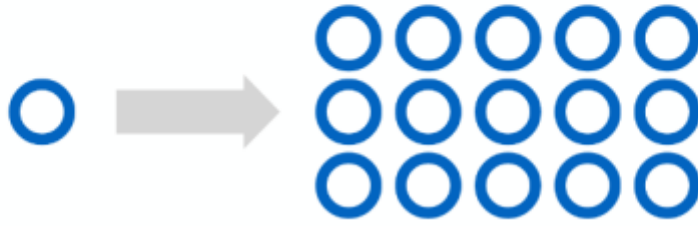
Problems that require you to extract a subset, sample, or summary from a Big Data source. You may do further analytics on the subset, and the subset might itself be quite large.



© 2015 RStudio, Inc. All rights reserved.

## Class 2. Compute on the parts

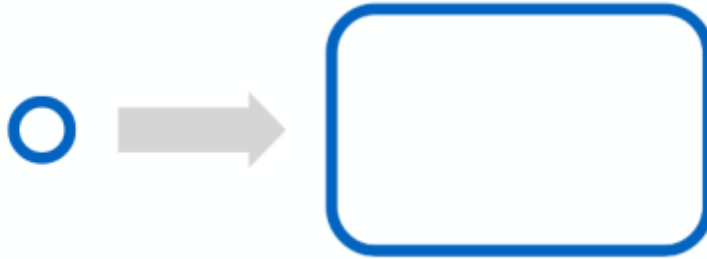
Problems that require you to repeat computation for many subgroups of the data, e.g. you need to fit one model per individual for thousands of individuals. You may combine the results once finished.



© 2015 RStudio, Inc. All rights reserved.

## Class 3. Compute on the whole

Problems that require you to use all of the data at once. These problems are irretrievably big; they must be run at scale within the data warehouse.



© 2015 RStudio, Inc. All rights reserved.

From an RStudio seminar<sup>[46]</sup> by Garrett Grolemond<sup>[47]</sup>

R interacts with other technologies





© 2015 RStudio, Inc. All rights reserved.

From an RStudio seminar<sup>[48]</sup> by Garrett Grolemund<sup>[49]</sup>

## R General Strategy (i): Class 1 & 2

### General Strategy

Store big data in a data warehouse

1. Pass subsets of data from warehouse to R
2. Transform R code, pass to warehouse.



© 2015 RStudio, Inc. All rights reserved.

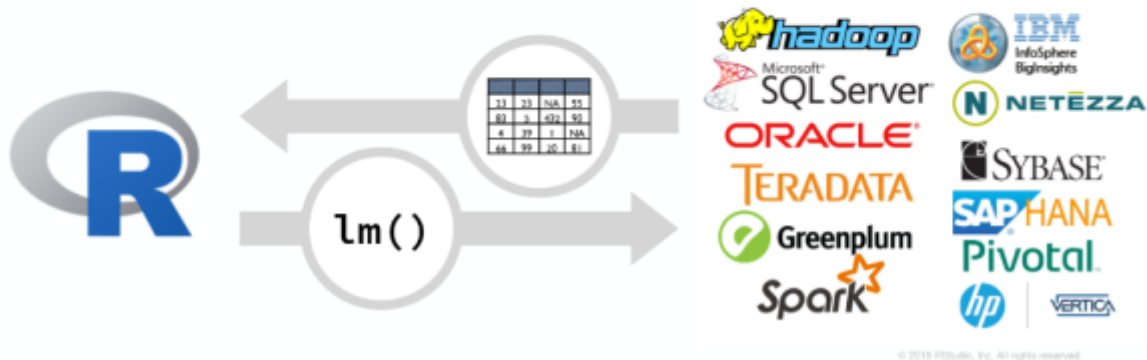
From an RStudio seminar<sup>[50]</sup> by Garrett Grolemund<sup>[51]</sup>

## R General Strategy (ii): Class 1 & 2

### General Strategy

Store big data in a data warehouse

1. Pass subsets of data from warehouse to R
2. Transform R code, pass to warehouse.



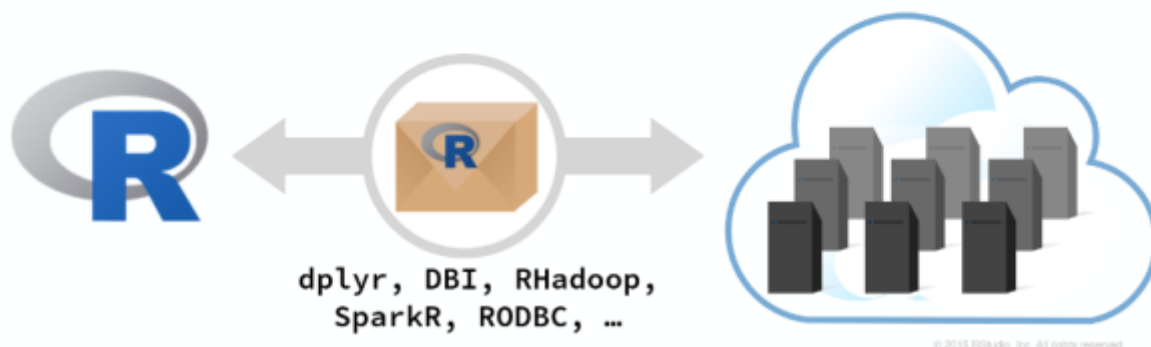
From an RStudio seminar<sup>[52]</sup> by Garrett Grolemund<sup>[53]</sup>

## R General Strategy (iii): Class 3

### General Strategy

Store big data in a data warehouse

1. Pass subsets of data from warehouse to R
2. Transform R code, pass to warehouse.



From an RStudio seminar<sup>[54]</sup> by Garrett Grolemund<sup>[55]</sup>

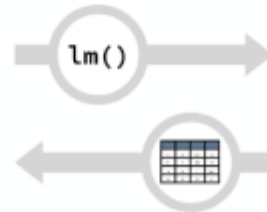
## Class 3 - dbplyr

R - tidyverse - dplyr + dbplyr



Package that provides data manipulation syntax for R. Comes with **built-in SQL backend**:

1. **Connects** to DBMS's
2. **Transforms R code** to SQL, sends to DBMS to run in DBMS
3. **Collect results** into R



© 2015 RStudio, Inc. All rights reserved.

From an RStudio seminar<sup>[56]</sup> by Garrett Grolemund<sup>[57]</sup>

## R General Strategy (iv): Recap

### Recap: Access Big Data



Store data outside of memory in data warehouse



Use an R package as an API to the data warehouse. dplyr, DBI, sparkR, others.

```
db <-
```

Create and work with connection object

```
rm(db)  
gc()
```

Close connection when finished

© 2015 RStudio, Inc. All rights reserved.

From an RStudio seminar<sup>[58]</sup> by Garrett Grolemund<sup>[59]</sup>

## (5) Spark (Apache) & sparklyr (Rstudio)

- Spark<sup>[60]</sup>, from Apache Foundation<sup>[61]</sup>: a leading tool that is democratizing our ability to process large datasets.
- R Packages:
  - Apache introduced an interface for the R computing language: **SparkR** R package
  - RStudio introduced **sparklyr** R package: a project merging R and Spark into a powerful

tool that is easily accessible to all.

- But **why where they introduced?**
- And **why 2 different packages?**

Derived from [here](#)<sup>[62]</sup> & [here](#)<sup>[63]</sup>

## Digital information vs analog information | World Bank - 2003

- World Bank report: digital information surpassed analog information around 2003.
  - 10 million terabytes of digital information (~ 10 million storage drives today)
  - our footprint of digital information is growing at exponential rates.



From "The R In Spark" (book)<sup>[64]</sup>

## Google File System | Google - 2003

- Search engines were unable to store all of the web page information required to support web searches in a single computer.
- They had to split information into several files and store them across many machines
- This approach became known as "**Google File System**", due to a research paper published in 2003 by Google.

From "The R In Spark" (book)<sup>[65]</sup>

## MapReduce | Google - 2004

- 2004: Google published a new paper describing how to perform operations across the Google File System:
  - approach called "MapReduce"
    - map operation: arbitrary way to transform each file into a new file
    - reduce operation: combines two files
  - Both operations require custom computer code, but the MapReduce framework takes care of automatically executing them across many computers at once.
  - Sufficient to process all the data available on the web, while providing flexibility to extract meaningful information from it.

From "The R In Spark" (book)<sup>[66]</sup>

# Hadoop Distributed File System (HDFS) | Yahoo - 2006

- A team at Yahoo implemented the Google File System and MapReduce as a single open source project, released in 2006 as **Hadoop**
  - Google File System implemented as the Hadoop Distributed File System (**HDFS**).
- The Hadoop project made distributed file-based computing accessible to a wider range of users and organizations, making MapReduce useful beyond web data processing.



From "The R In Spark" (book)<sup>[67]</sup> | Image from De Apache Software Foundation, with Apache License 2.0<sup>[68]</sup>

# Hive | Facebook - 2008

- Hadoop provided support to perform MapReduce operations over a distributed file system, but it still required MapReduce operations to be written with code every time a data analysis was run
- **Hive** project (2008, by Facebook) brought **Structured Query Language** (SQL) support to Hadoop.
  - Data analysis could now be performed at large scale without the need to write code for each MapReduce operation



From "The R In Spark" (book)<sup>[69]</sup> | Image from Davod<sup>[70]</sup> with Apache License 2.0<sup>[71]</sup>

# Spark (closed sourced) | UC Berkely - 2009

- In 2009, **Spark** began as a research project at UC Berkeley's AMPLab to improve on MapReduce.
  - Spark provided a richer set of verbs beyond MapReduce to facilitate optimizing code

running in multiple machines.

- **Spark also loaded data in-memory**, making operations much faster than **Hadoop's on-disk storage**.



From "The R In Spark" (book)<sup>[72]</sup>

## Spark: In-memory and on-disk

- Spark: well known for its in-memory performance, but designed to be a general execution engine that works both in-memory and on-disk
  - For instance, Spark has set sorting, for which data was not loaded in-memory, but using improvements in network serialization, network shuffling, and efficient use of the CPU's cache to dramatically enhance performance.
  - If you needed to sort large amounts of data, there was no other system in the world faster than Spark.

From "The R In Spark" (book)<sup>[73]</sup>

## Spark: faster and easier

- Spark is much faster, more efficient, and easier to use than Hadoop.
- Speed Example:
  - Without Spark: it takes 72 minutes and 2,100 computers to sort 100 terabytes of data using Hadoop,
  - With Spark: only 23 minutes and 206 computers
- Simplicity example: word-counting MapReduce example takes:
  - about 50 lines of code in Hadoop, but
  - only 2 lines of code in Spark.

	Hadoop Record	Spark Record
Data Size	102.5 TB	100 TB
Elapsed Time	72 mins	23 mins
Nodes	2100	206
Cores	50400	6592
Disk	3150 GB/s	618 GB/s
Network	10Gbps	10Gbps
Sort rate	1.42 TB/min	4.27 TB/min
Sort rate / node	0.67 GB/min	20.7 GB/min

From "The R In Spark" (book)<sup>[74]</sup>



# Spark (open sourced) - 2010 | Apache Foundation - 2013

- 2010: open sourced. 2013: donated to the Apache Software Foundation
- Apache Spark is a unified analytics engine for large-scale data processing
  - Unified: supports many libraries, cluster technologies, and storage systems.
  - Analytics: discovery & interpretation of data to communicate information
  - Engine: expected to be efficient and generic.
  - Large-Scale: as cluster-scale (set of connected computers working together).



From "The R In Spark" (book)<sup>[75]</sup> | Image from the Apache Software Foundation, with Apache License 2.0<sup>[76]</sup>

## SparkR (Base R) vs. sparklyr (Modern R)

Feature	SparkR	sparklyr
Data input & output	++	++
Data manipulation	-	+++
Documentation	++	++
Ease of setup	++	++
Function naming	--	+++
Installation	+	++
Machine learning	+	++
Range of functions	+++	++
Running arbitrary code	+	++
Tidyverse compatability	---	+++

From <https://www.eddjberry.com/post/2017-12-05-sparkr-vs-sparklyr/><sup>[77]</sup>

## Sparklyr | RStudio

- Sparklyr, from Rstudio: <https://spark.rstudio.com/><sup>[78]</sup>
- R interface for Apache Spark, agnostic to Spark versions,
  - 2016: 1st version released (v0.4)
- Easy to install, serving the R community,
- Embracing other packages & practices from the R community (tidyverse, ...)
- Designed for: New Users, Data Scientists, and Expert Users



Image from here<sup>[79]</sup>

## Summary: Big Data with Modern R & Spark

Big Data with Modern R & Spark in context

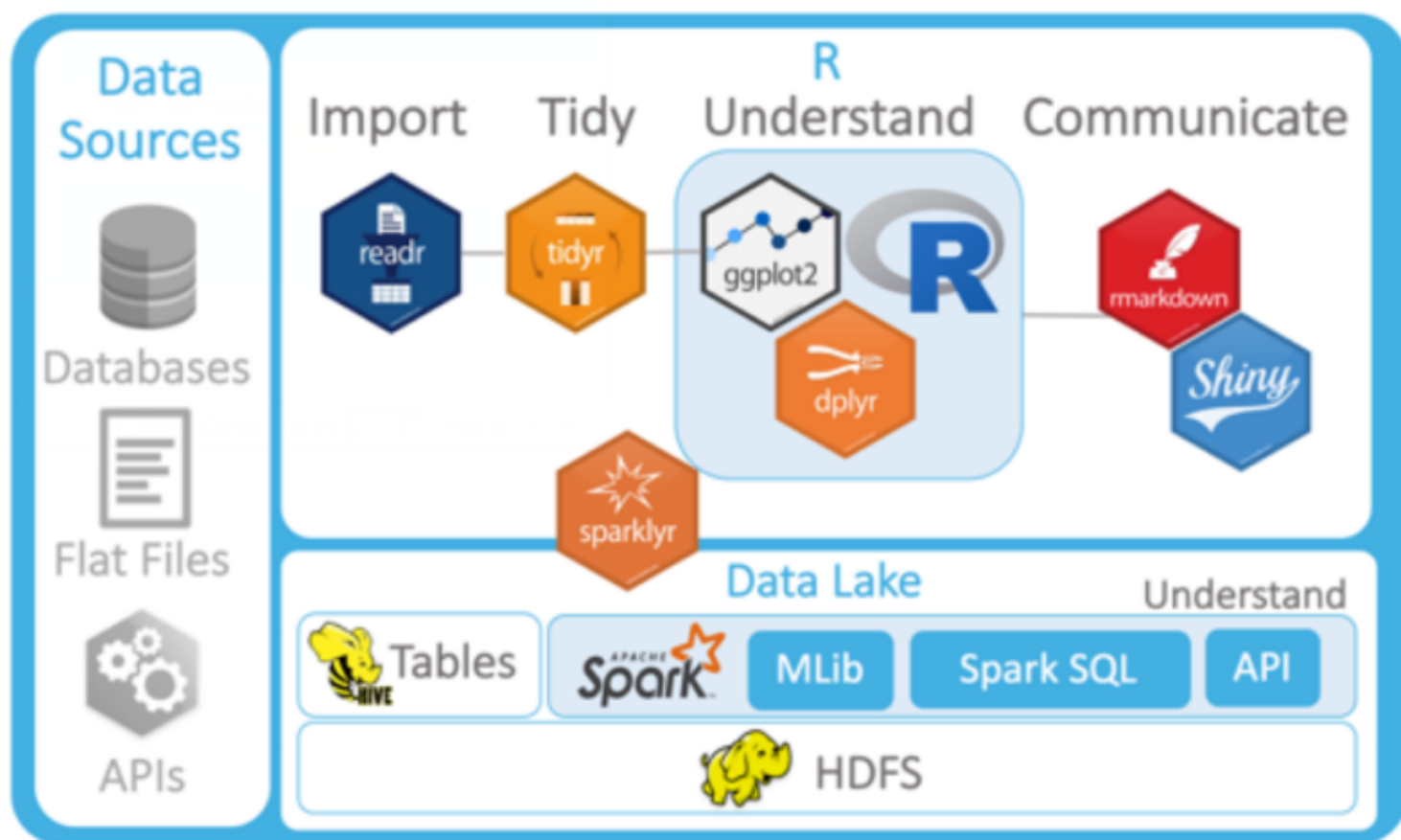


Image from here<sup>[80]</sup>

## Speeding up Spark via R with Apache Arrow

Apache Arrow<sup>[81]</sup> is a cross-language development platform for in-memory data, you can read more about this in the Arrow and beyond<sup>[82]</sup> blog post. In sparklyr 1.0, we are embracing Arrow as an efficient bridge between R and Spark, conceptually.



X Axis: Time to complete task (left-hand-side is faster)

Y Axis: Top section: **WITHOUT** Apache Arrow), Bottom section: **WITH** apache Arrow.

Copying:	Collecting:	Transforming
		

From: <https://arrow.apache.org/blog/2019/01/25/r-spark-improvements/><sup>[83]</sup>


## (6) Sparklyr - next steps



From: <https://github.com/harryprince/awesome-sparklyr><sup>[84]</sup>

## (7) Addendum: Diskframe as potential alternative for medium sized projects?

- `disk.frame` package<sup>[85]</sup>, answer to: how do I manipulate structured tabular data that doesn't fit into Random Access Memory (RAM)?.

- It makes use of two simple ideas  **disk.frame**  
FOR >> RAM DATA

1. split up a larger-than-RAM dataset into chunks and store each chunk in a separate file inside a folder and
  2. provide a convenient API to manipulate these chunks
- It performs a similar role to distributed systems such as Apache Spark, Python's Dask, and Julia's JuliaDB.jl for medium data which are datasets that are too large for RAM but not quite large

enough to qualify as *big data*.

More information: [here](#)<sup>[86]</sup> and [here](#)<sup>[87]</sup>

## (8) Former hands-on exercise nowadays with sparklyr

Remember the hands on exercise we did on the session about "**Reproducible Work in Data Science**" in this postgraduate course?

We will revisit the work we did in that session, and we will evolve the Rmd you did in order to adapt it to use sparklyr to achieve the goal we had in that hands-on exercise.

We can't make use of the posit.cloud free plan since the Spark infrastructure we need does require more than just 1 Gb of RAM (the RAM provided for free in posit.cloud free plan).

So that in this seminar, we will see the job performed in [datascience.seeds4c.org](https://datascience.seeds4c.org), which is a lxc (linux container) based on Ubuntu 20.04 LTS with 4Gb of RAM and 4 cpu.

<http://datascience.seeds4c.org:8787/><sup>[88]</sup>

We will run this code chunk by chunk, to assess the RAM and time consumed to perform the same task with and without sparklyr in the same container.

[https://gitlab.com/xavidp/datascience2023/-/blob/master/DATA\\_SMC\\_with\\_sparklyr.Rmd](https://gitlab.com/xavidp/datascience2023/-/blob/master/DATA_SMC_with_sparklyr.Rmd)<sup>[89]</sup>

Task: just loading the whole dataset (600 Mb csv file on disk), for instance, and saving it partitioned by meteorological station to csv files on disk:

- From **117 Mb (and 63 secs)** in the R Session using **sparklyr** to **1560 Mb (and 165 secs)** if using **just R (not involving Spark)**.

## References

- TheRinSpark Book: <https://therinspark.com/><sup>[90]</sup>
- Cheatsheets at <https://www.rstudio.com/resources/cheatsheets/><sup>[91]</sup>
  - And some in Spanish Cheatsheets URL [URL](#)<sup>[92]</sup> > Spanish Translations - Traducciones en español
- Posit (Rstudio) Cloud (with free plan): <https://posit.cloud/plans><sup>[93]</sup>
- Webinars & scripts for hands-on practising:
  - <https://spark.rstudio.com/get-started/><sup>[94]</sup>
  - <https://gitlab.com/radup/curs-r-introduccio/blob/master/codi/extra.tips.bigdata.R><sup>[95]</sup>
  - <https://diskframe.com/articles/02-intro-disk-frame.html><sup>[96]</sup>
- Compilation of links related to sparklyr: <https://github.com/harryprince/awesome-sparklyr><sup>[97]</sup>

- Playground for R & Spark (14d for free): <https://community.cloud.databricks.com/><sup>[98]</sup> Pricing<sup>[99]</sup>
- Tidyverse Skills for Data Science. Carrie Wright, Shannon E. Ellis, Stephanie C. Hicks and Roger D. Peng. 2021-09-02
  - <https://jhudatascience.org/tidyversecourse/><sup>[100]</sup>

# Thanks

<sup>[101]</sup>

[xavier.depedro \(a\) seeds4c.org](https://xavier.depedro(a)seeds4c.org)

Unless elsewhere noted, contents of this web site are released under a Creative Commons<sup>[102]</sup> license.

## Earlier videorecording

[+]

---

<sup>[1]</sup> <https://www.ub.edu>

<sup>[2]</sup> <https://www.ub.edu>

<sup>[3]</sup> <https://www.uoc.edu>

<sup>[4]</sup> <https://www.barcelona.cat/en/>

<sup>[5]</sup> <https://ajuntament.barcelona.cat/imi/en>

<sup>[6]</sup> <https://www.barcelona.cat/en/>

<sup>[7]</sup> <https://ajuntament.barcelona.cat/digital/en/digital-transformation/city-data-commons/municipal-data-office>

<sup>[8]</sup> <https://www.barcelona.cat/en/>

<sup>[9]</sup> <http://ueb.vhir.org/>

<sup>[10]</sup> <https://en.vhir.org>

<sup>[11]</sup> <http://ueb.vhir.org/>

<sup>[12]</sup> <https://en.vhir.org>

<sup>[13]</sup> <https://ajuntament.barcelona.cat/digital/en/digital-transformation/city-data-commons/municipal-data-office>

<sup>[14]</sup> <https://ajuntament.barcelona.cat/premsa/2018/02/13/barcelona-crea-una-nova-oficina-municipal-de-dades/>

<sup>[15]</sup> <https://www.bcn.cat/estadistica/angles/>

<sup>[16]</sup> [https://ajuntament.barcelona.cat/digital/en/noticia/a-website-has-been-launched-which-carries-out-analytic-monitoring-of-the-evolution-of-covid-19-in-the-city\\_956022](https://ajuntament.barcelona.cat/digital/en/noticia/a-website-has-been-launched-which-carries-out-analytic-monitoring-of-the-evolution-of-covid-19-in-the-city_956022)

<sup>[17]</sup> <https://dades.ajuntament.barcelona.cat/seguiment-covid19-bcn/>

<sup>[18]</sup> <https://publiccode.eu/>

<sup>[19]</sup> <https://ajuntament.barcelona.cat/digital/en/digital-transformation/technology-for-a-better-government/open-source-software>

<sup>[20]</sup> [http://ajuntament.barcelona.cat/digital/sites/default/files/LE\\_MesuradeGovern\\_EN\\_9en.pdf](http://ajuntament.barcelona.cat/digital/sites/default/files/LE_MesuradeGovern_EN_9en.pdf)

<sup>[21]</sup> <https://ajuntament.barcelona.cat/digital/en/digital-transformation/city-data-commons/cityos>

<sup>[22]</sup> <https://www.youtube.com/watch?v=OsbpZTzE5SI>

<sup>[23]</sup> <https://www.slideshare.net/francescabria/bria-francesca-bcn-open-source-agile-digital-transformation-strategy>

- [24] [https://github.com/AjuntamentdeBarcelona/CityOS\\_AjBCN](https://github.com/AjuntamentdeBarcelona/CityOS_AjBCN)
- [25] [https://github.com/AjuntamentdeBarcelona/CityOS\\_AjBCN/blob/master/doc/COS1%20-%20SP12\\_Disseny%20de%20Seguretat%20del%20sistema%20inform%C3%A0tic.docx](https://github.com/AjuntamentdeBarcelona/CityOS_AjBCN/blob/master/doc/COS1%20-%20SP12_Disseny%20de%20Seguretat%20del%20sistema%20inform%C3%A0tic.docx)
- [26] [https://github.com/AjuntamentdeBarcelona/CityOS\\_AjBCN](https://github.com/AjuntamentdeBarcelona/CityOS_AjBCN)
- [27] [https://github.com/AjuntamentdeBarcelona/CityOS\\_AjBCN/blob/master/doc/COS1%20-%20SP12\\_Disseny%20de%20Seguretat%20del%20sistema%20inform%C3%A0tic.docx](https://github.com/AjuntamentdeBarcelona/CityOS_AjBCN/blob/master/doc/COS1%20-%20SP12_Disseny%20de%20Seguretat%20del%20sistema%20inform%C3%A0tic.docx)
- [28] <https://cran.r-project.org/web/packages/tidyverse/vignettes/paper.html>
- [29] <https://cran.r-project.org/src/contrib/Archive/tidyverse/>
- [30] <http://r4stats.com/workshops/introduction-to-modern-r/>
- [31] <https://www.infoq.com/podcasts/martin-hadley-r-ecosystem>
- [32] <https://www.amazon.es/Modern-Science-Chapman-Texts-Statistical-ebook/dp/B06XPNZ4H2>
- [33] [https://b-rodrigues.github.io/modern\\_R/](https://b-rodrigues.github.io/modern_R/)
- [34] <https://www.datascienceworkshops.com/catalogue/modern-r-and-the-tidyverse/>
- [35] <https://www.youtube.com/watch?v=a-Yb518580c>
- [36] <https://moderndive.com/>
- [37] <http://modern-rstats.eu/>
- [38] <https://rviews.rstudio.com/2017/06/08/what-is-the-tidyverse/>
- [39] <https://github.com/rstudio/webinars/tree/master/55-ciencia-de-datos-R>
- [40] <https://www.r-bloggers.com/2021/05/the-new-r-pipe/>
- [41] <https://cran.r-project.org/web/packages/tidyverse/vignettes/paper.html>
- [42] <https://raw.githubusercontent.com/rstudio/webinars/master/05-Data-Wrangling-with-R-and-RStudio/wrangling-webinar.pdf>
- [43] <https://raw.githubusercontent.com/rstudio/webinars/master/55-ciencia-de-datos-R/ciencia-de-datos-R.pdf>
- [44] <https://github.com/rstudio/webinars/blob/master/14-Work-with-big-data/14-Work-with-big-data.pdf>
- [45] <https://github.com/garrettgman>
- [46] <https://github.com/rstudio/webinars/blob/master/14-Work-with-big-data/14-Work-with-big-data.pdf>
- [47] <https://github.com/garrettgman>
- [48] <https://github.com/rstudio/webinars/blob/master/14-Work-with-big-data/14-Work-with-big-data.pdf>
- [49] <https://github.com/garrettgman>
- [50] <https://github.com/rstudio/webinars/blob/master/14-Work-with-big-data/14-Work-with-big-data.pdf>
- [51] <https://github.com/garrettgman>
- [52] <https://github.com/rstudio/webinars/blob/master/14-Work-with-big-data/14-Work-with-big-data.pdf>
- [53] <https://github.com/garrettgman>
- [54] <https://github.com/rstudio/webinars/blob/master/14-Work-with-big-data/14-Work-with-big-data.pdf>
- [55] <https://github.com/garrettgman>
- [56] <https://github.com/rstudio/webinars/blob/master/14-Work-with-big-data/14-Work-with-big-data.pdf>
- [57] <https://github.com/garrettgman>
- [58] <https://github.com/rstudio/webinars/blob/master/14-Work-with-big-data/14-Work-with-big-data.pdf>
- [59] <https://github.com/garrettgman>
- [60] <https://spark.apache.org/>
- [61] <https://apache.org/>
- [62] <https://github.com/rstudio/webinars/blob/master/42-Introduction%20to%20sparklyr/Introducing%20sparklyr%20-%20Webinar.pdf>



<sup>[63]</sup> <https://therinspark.com/>

<sup>[64]</sup> <https://therinspark.com/>

<sup>[65]</sup> <https://therinspark.com/>

<sup>[66]</sup> <https://therinspark.com/>

<sup>[67]</sup> <https://therinspark.com/>

<sup>[68]</sup> [https://svn.apache.org/repos/asf/hadoop/logos/out\\_rgb/](https://svn.apache.org/repos/asf/hadoop/logos/out_rgb/)

<sup>[69]</sup> <https://therinspark.com/>

<sup>[70]</sup> [https://commons.wikimedia.org/wiki/User:Amitie\\_10g](https://commons.wikimedia.org/wiki/User:Amitie_10g)

<sup>[71]</sup> <http://www.apache.org/licenses/LICENSE-2.0>

<sup>[72]</sup> <https://therinspark.com/>

<sup>[73]</sup> <https://therinspark.com/>

<sup>[74]</sup> <https://therinspark.com/>

<sup>[75]</sup> <https://therinspark.com/>

<sup>[76]</sup> [http://svn.apache.org/repos/asf/hadoop/logos/asf\\_hadoop/](http://svn.apache.org/repos/asf/hadoop/logos/asf_hadoop/)

<sup>[77]</sup> <https://www.eddjberrry.com/post/2017-12-05-sparkr-vs-sparklyr/>

<sup>[78]</sup> <https://spark.rstudio.com/>

<sup>[79]</sup> <https://www.slideshare.net/ICTeam/sparklyr-big-data-enabler-for-r-users>

<sup>[80]</sup> <https://www.slideshare.net/ICTeam/sparklyr-big-data-enabler-for-r-users>

<sup>[81]</sup> <https://arrow.apache.org/>

<sup>[82]</sup> <https://blog.rstudio.com/2018/04/19/arrow-and-beyond/>

<sup>[83]</sup> <https://arrow.apache.org/blog/2019/01/25/r-spark-improvements/>

<sup>[84]</sup> <https://github.com/harryprince/awesome-sparklyr>

<sup>[85]</sup> <https://github.com/xiaodaigh/disk.frame>

<sup>[86]</sup> <https://github.com/xiaodaigh/disk.frame>

<sup>[87]</sup> <https://diskframe.com/articles/intro-disk-frame.html>

<sup>[88]</sup> <http://datascience.seeds4c.org:8787/>

<sup>[89]</sup> [https://gitlab.com/xavidp/datascience2023/-/blob/master/DATA\\_SMC\\_with\\_sparklyr.Rmd](https://gitlab.com/xavidp/datascience2023/-/blob/master/DATA_SMC_with_sparklyr.Rmd)

<sup>[90]</sup> <https://therinspark.com/>

<sup>[91]</sup> <https://www.rstudio.com/resources/cheatsheets/>

<sup>[92]</sup> <https://www.rstudio.com/resources/cheatsheets/>

<sup>[93]</sup> <https://posit.cloud/plans>

<sup>[94]</sup> <https://spark.rstudio.com/get-started/>

<sup>[95]</sup> <https://gitlab.com/radup/curs-r-introduccio/blob/master/codi/extra.tips.bigdata.R>

<sup>[96]</sup> <https://diskframe.com/articles/02-intro-disk-frame.html>

<sup>[97]</sup> <https://github.com/harryprince/awesome-sparklyr>

<sup>[98]</sup> <https://community.cloud.databricks.com/>

<sup>[99]</sup> <https://databricks.com/product/aws-pricing>

<sup>[100]</sup> <https://jhudatascience.org/tidyversecourse/>

<sup>[101]</sup> <http://creativecommons.org/licenses/by-sa/3.0/>

<sup>[102]</sup> <http://creativecommons.org/licenses/by-sa/3.0/>