

# R y Spark para la Ciencia de Datos

Edgar Ruiz

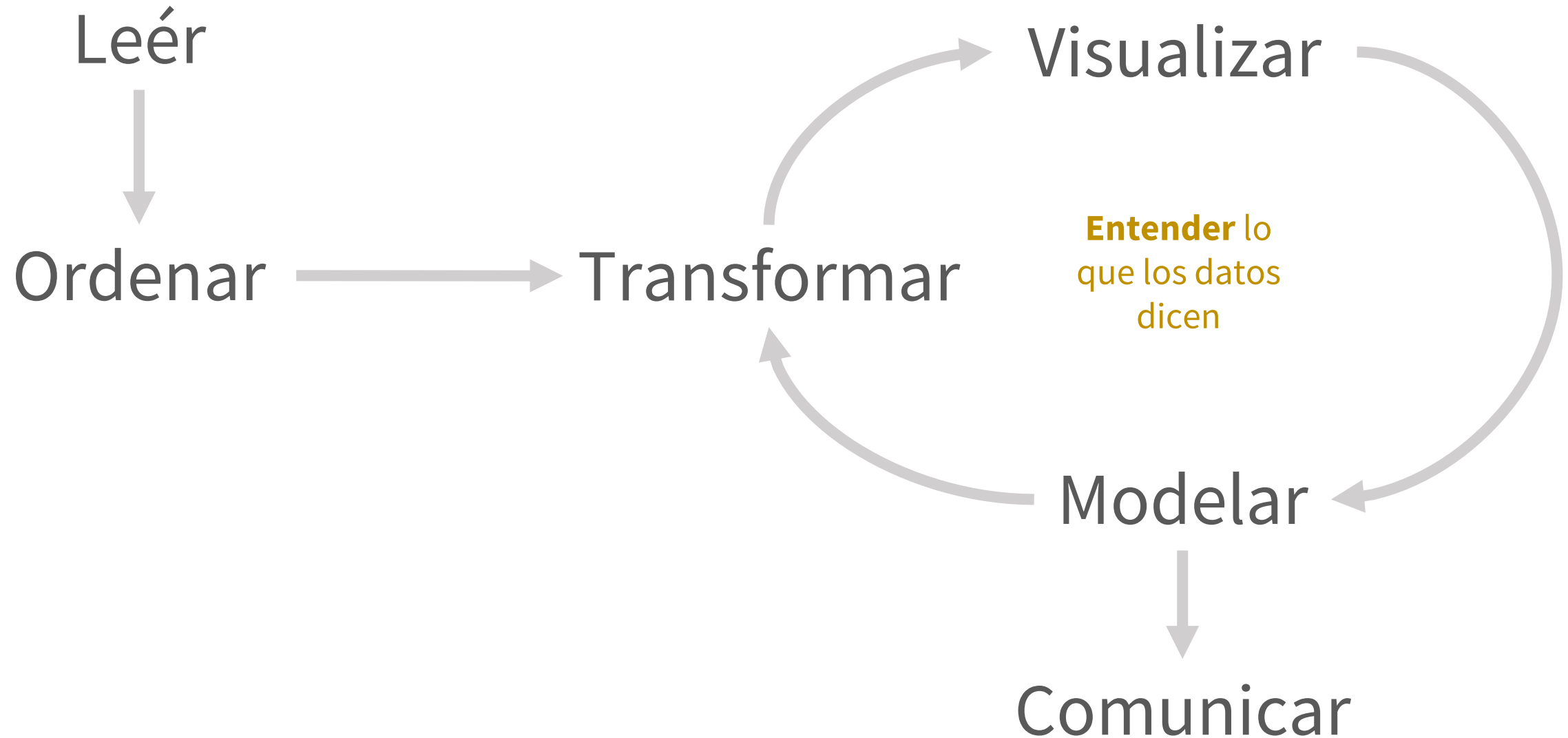
 @theotheredgar

 [linkedin.com/in/edgararuiz](https://www.linkedin.com/in/edgararuiz)

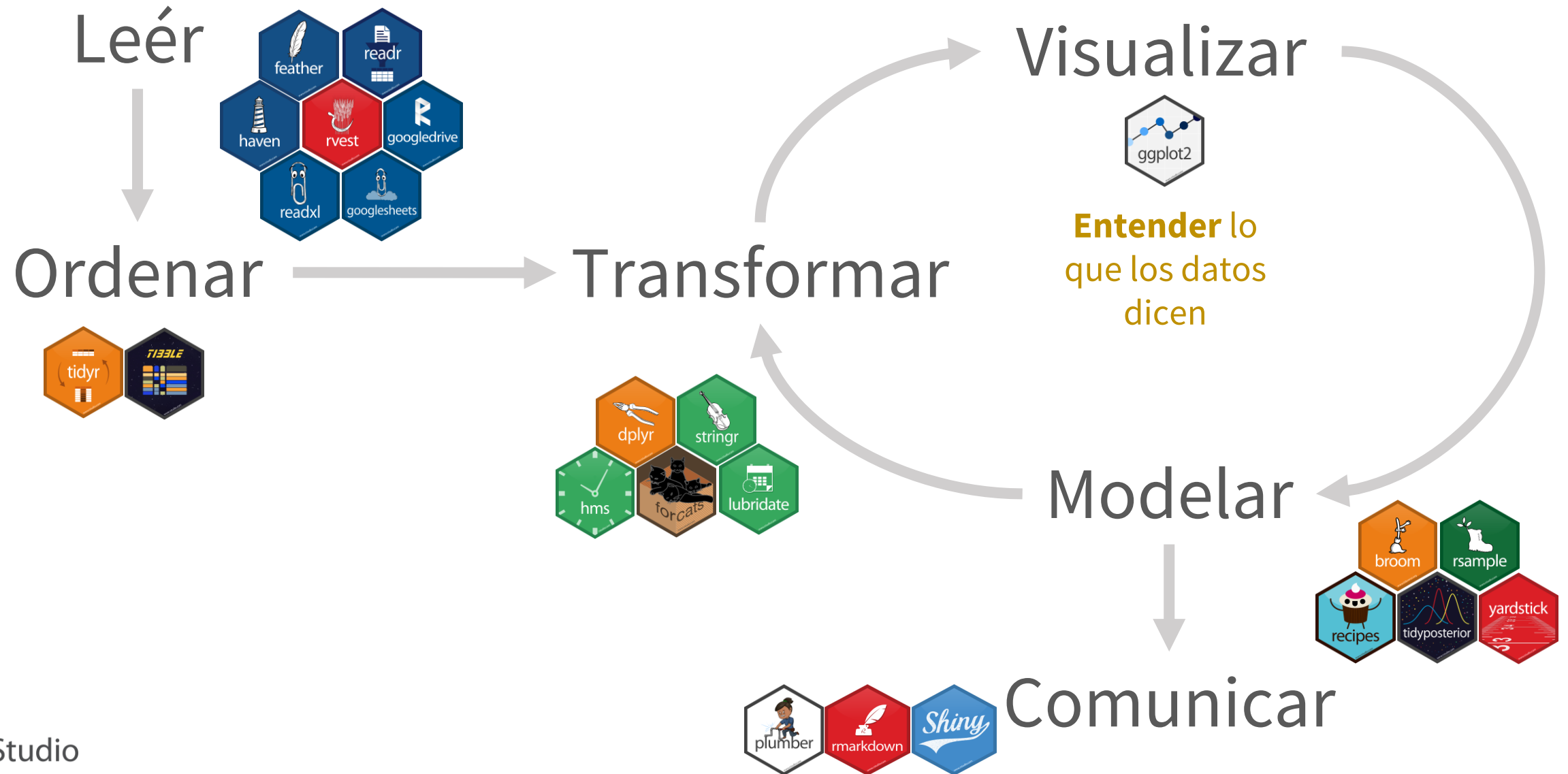
 [github.com/edgararuiz](https://github.com/edgararuiz)

24 de marzo del 2019

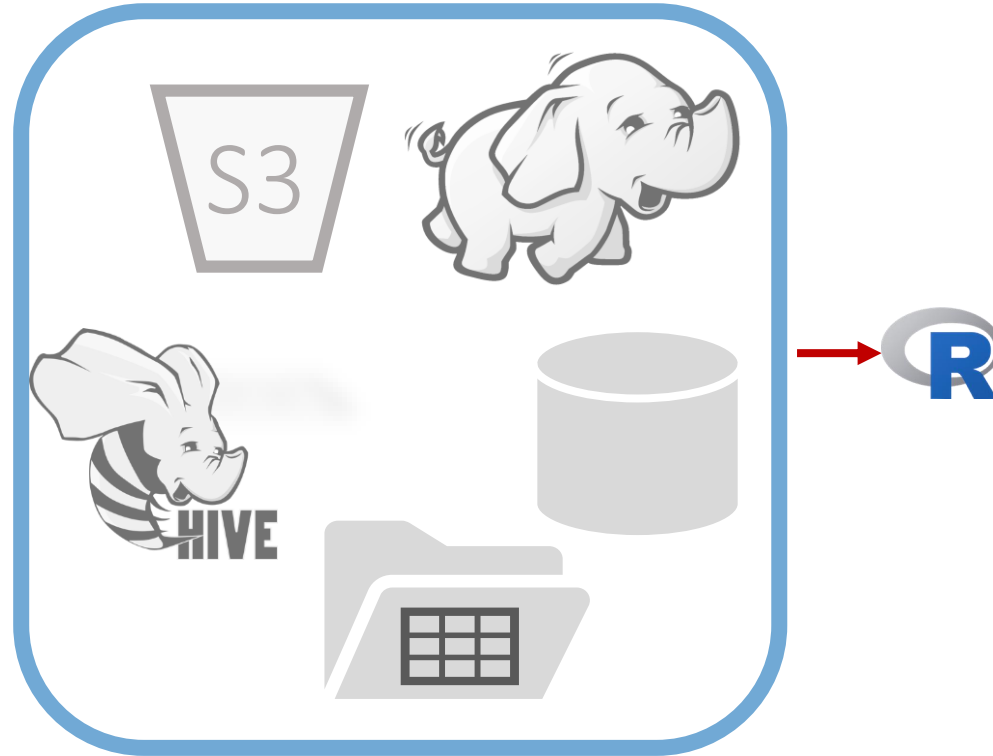
# Ciencia de Datos



# Todo se prepara dentro de



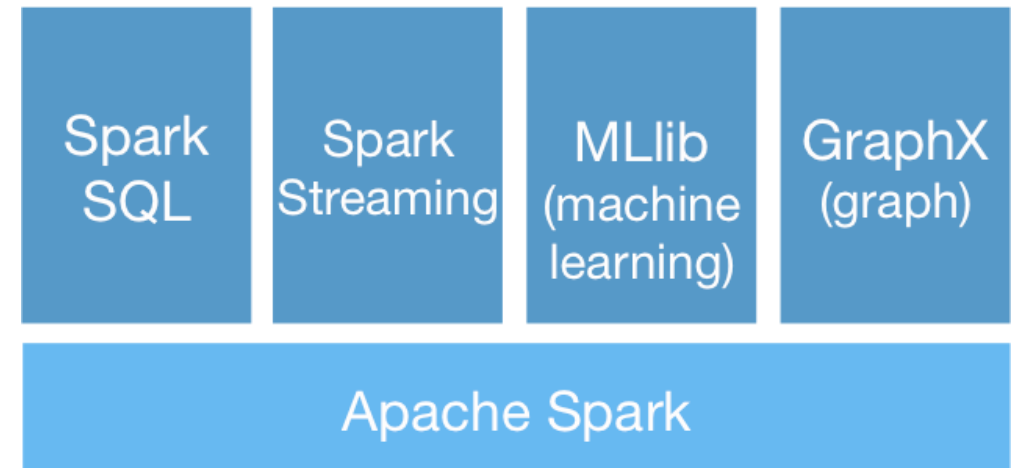
# Datos mas grandes que el RAM



# Que es ?

Motor analítico para procesar datos a gran escala

- Informática en clúster
- Aprendizaje automático
- Comunicación usando SQL
- API extensible

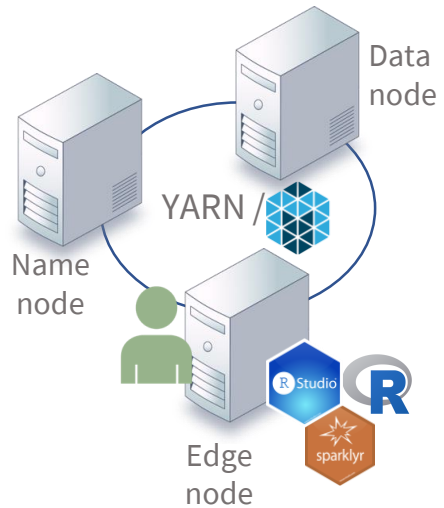


# Tipos de análisis disponibles en Spark

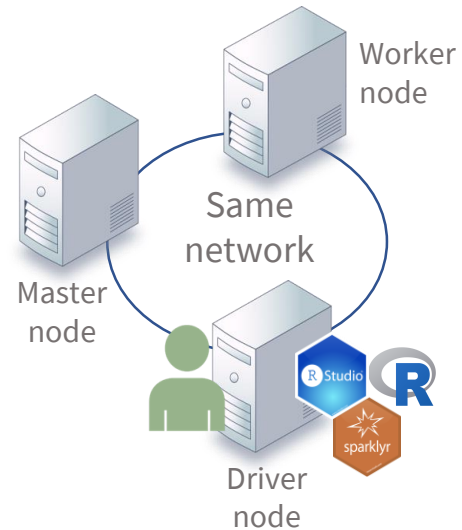
- Modelos de regresión
- Modelos de clasificación
- Modelos de agrupación (clustering)
- Modelos de gráficas (GraphX)
- Análisis sobre datos stream (constante flujo)
- Análisis de texto, incluyendo modelos

# Variedad de implementaciones

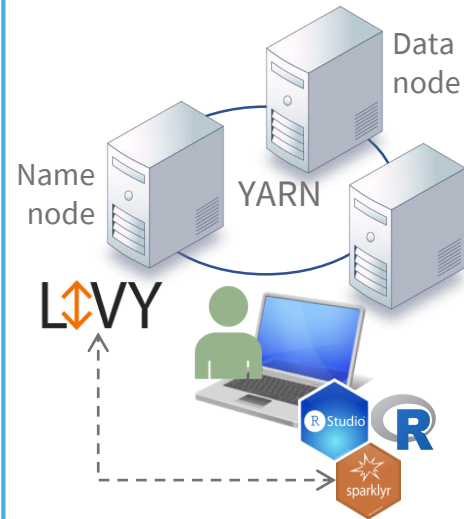
## Clusters manejados



## Clusters sólo con Spark



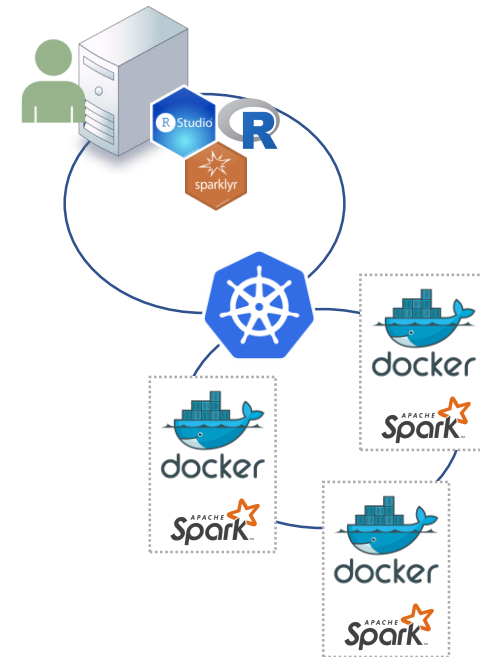
## Livy



## Local



## Kubernetes



# Que es ?

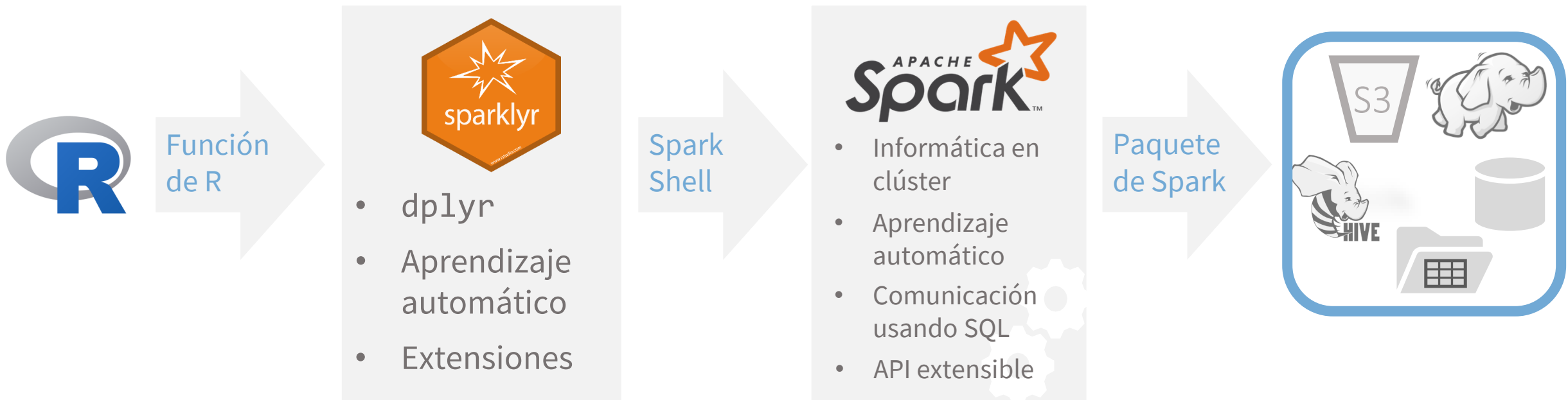
Es una interfaz para R y Spark

- Provee soporte a dplyr dentro de Spark
- Acceso a todo el API de Spark
- ...incluyendo Pipelines

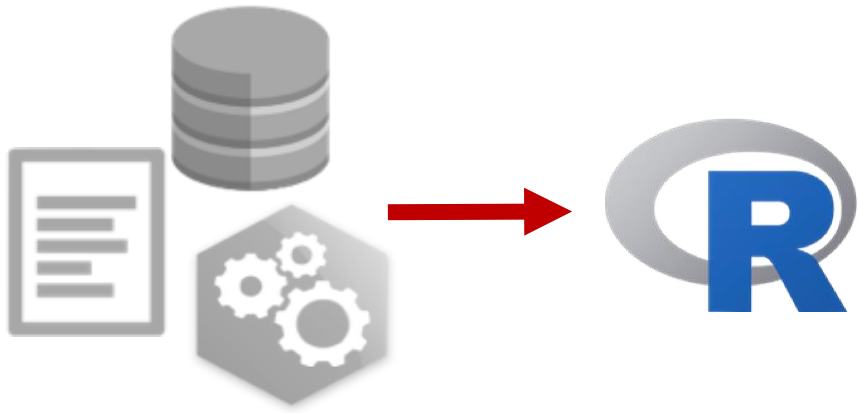




# Como funciona sparklyr



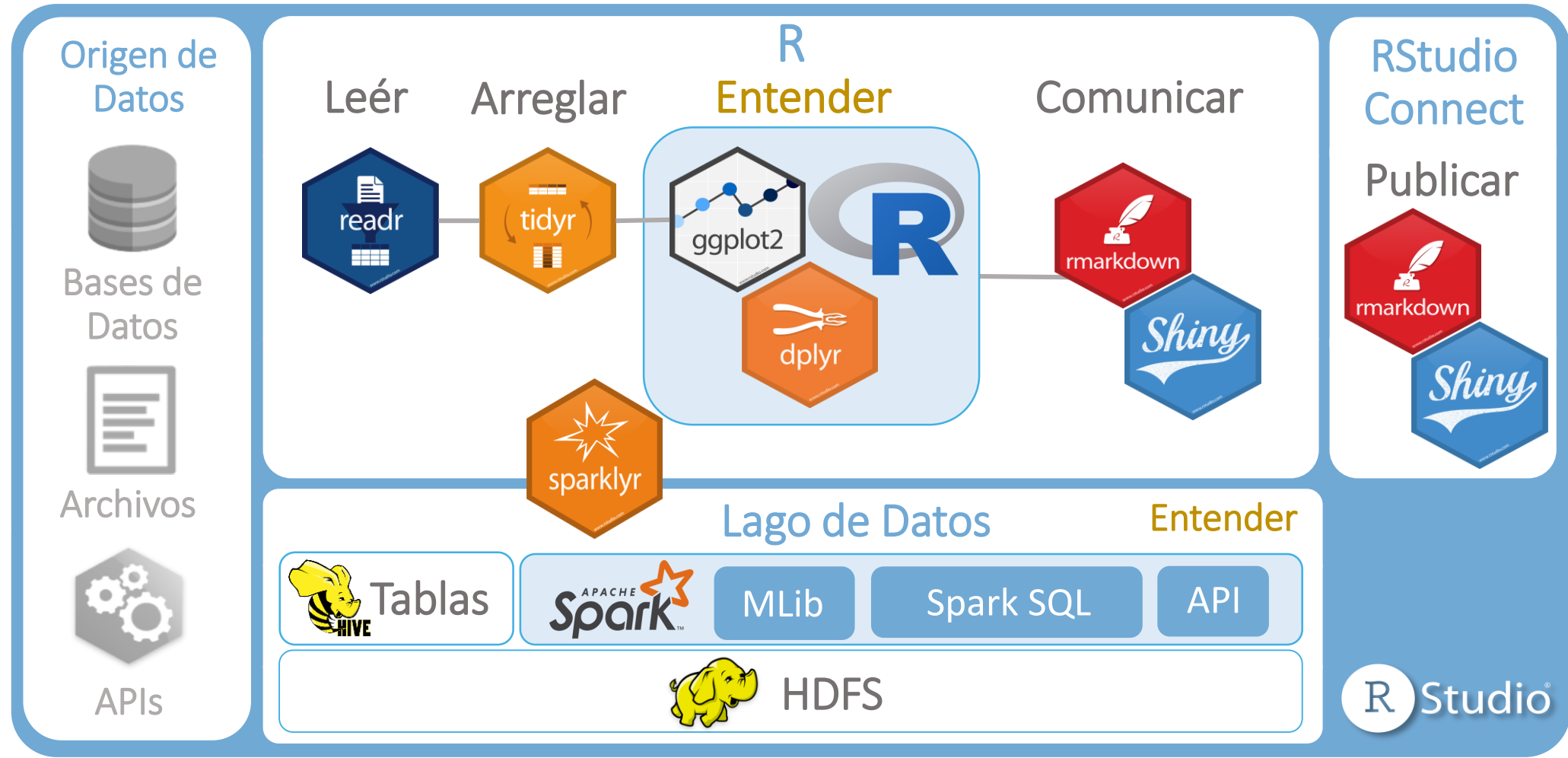
# La idea principal



Extraer  
datos



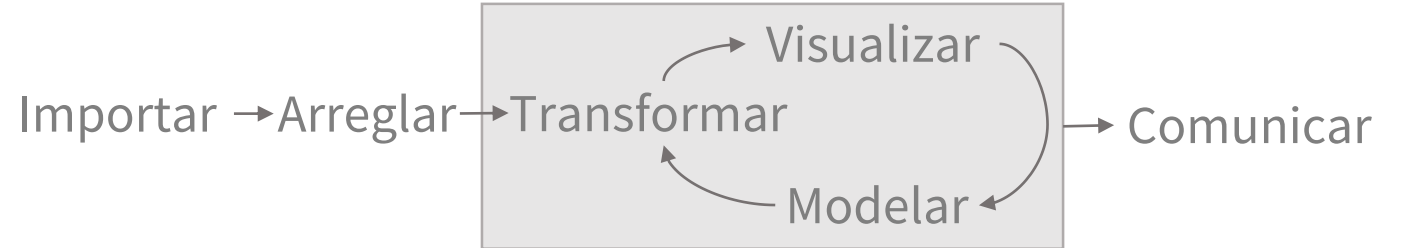
# Ciencia de Datos con R y Spark



# Mas allá que Ciencia de Datos...

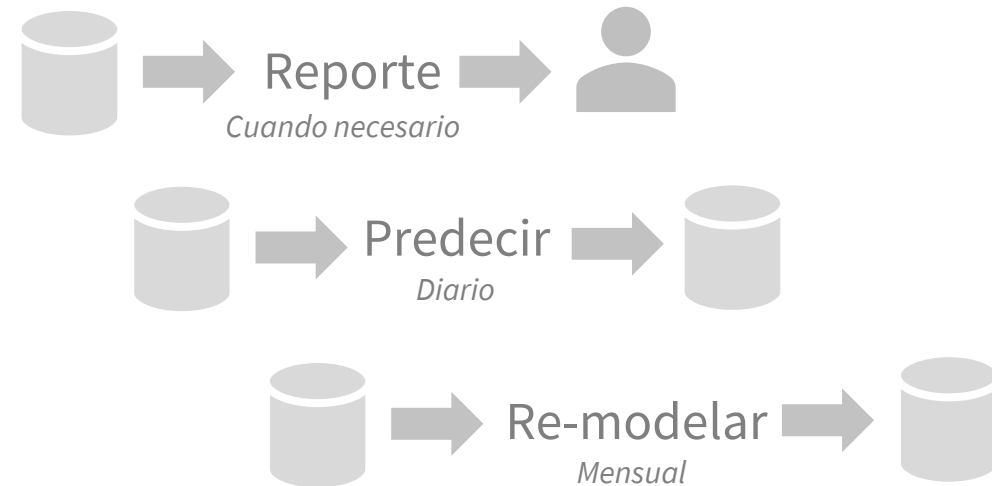
## Ciencia de Datos

- Resultado: Conocimiento
- Muchos experimentos
- Exploración



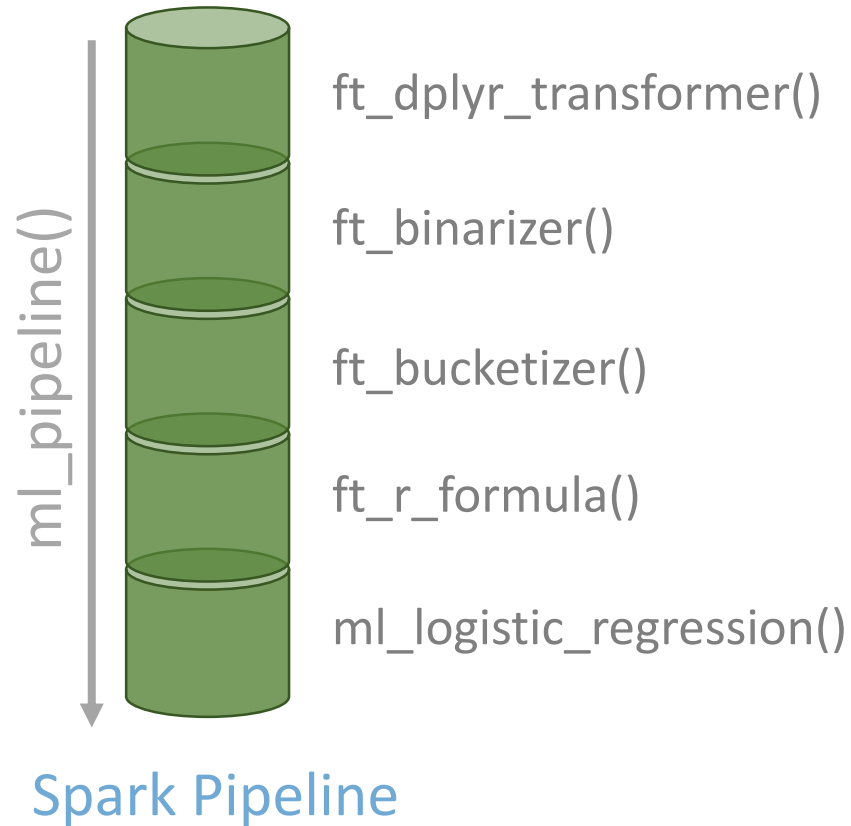
## Producción

- Resultado: **Software**
- Probado
- Automatizado
- Más controlado

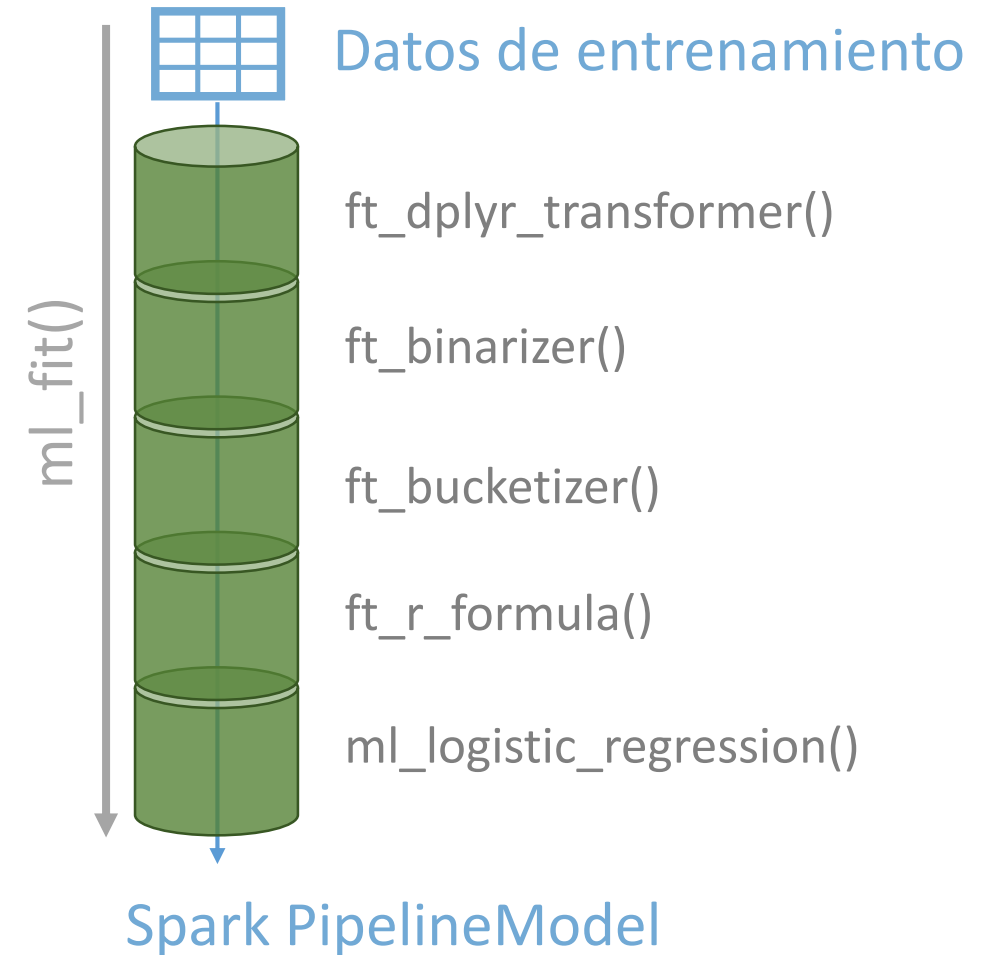


# Spark Pipelines (Tubería)

## Estimator



## Transformer



# Recursos

Sitio oficial de sparklyr:

[spark.rstudio.com](https://spark.rstudio.com)

Sitio oficial de Spark:

[spark.apache.org](https://spark.apache.org)



Para aprender como usar los paquetes en práctica, las Hojas de Referencia, o *Cheatsheets*, son los mejores recursos, no importa el idioma

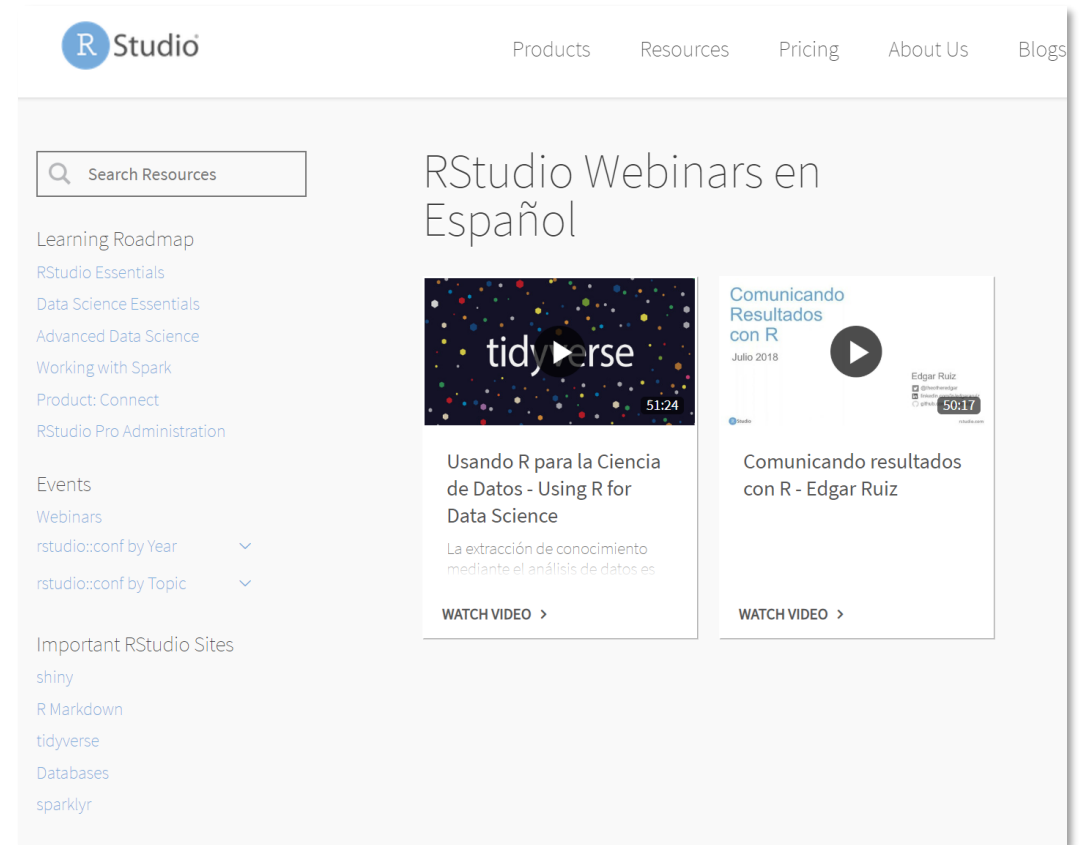
[rstudio.com/resources/cheatsheets](https://rstudio.com/resources/cheatsheets)



# Recursos en español - Webinars

Una base para aprender a utilizar R de manera efectiva en nuestros análisis. Estos webinars proveen tres cosas:

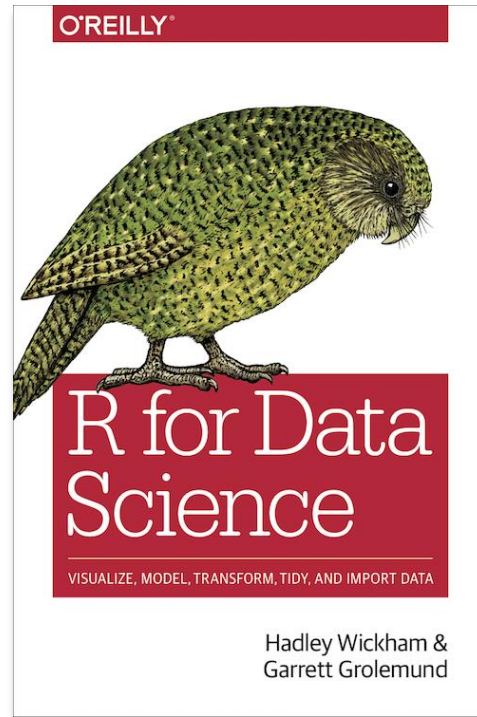
1. Ejemplos de código
2. Presentaciones
3. Video de la sesión



[resources.rstudio.com/espanol](https://resources.rstudio.com/espanol)



# Muy pronto!



[github.com/cienciadedatos](https://github.com/cienciadedatos)

# Materiales

[rstudio.io/conectar](https://rstudio.io/conectar)