# Big Data with R

## Edgar Ruiz

@theotheredgar
linkedin.com/in/edgararuiz

**May 2018**



RStudio

# Let's talk about **Big Data**

R Studio

Big data?

Velocity

Volume

Value

Variety

Veracity

Data > RAM

Garrett Grolemund

Remote Data

Edgar Ruiz

R Studio

db.rstudio.com / spark.rstudio.com

# Big Data in R



Data

RAM

Small Conduit

Remote Data

# Big Data Strategies

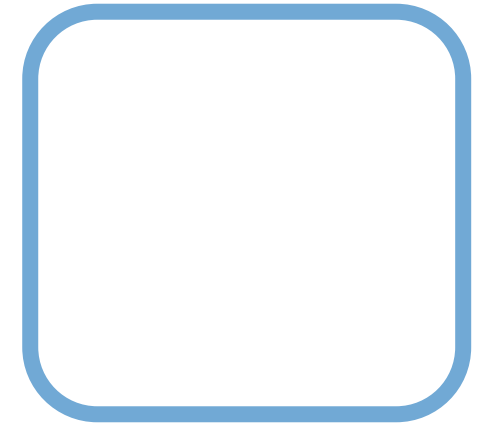| **Sample** | **Parts** | **Whole** |
|:---:|:---:|:---:|
| Most common approach for **modeling** | Most common approach for **general analysis** | In most cases, **the preferred approach,** it's just not feasible |

# Push compute, collect results

# Ideally, analyze in place

```
select count() from sales where amount > 1000 group by month
```

*"Number of sales over $1K by month"*

Returns a **data.frame** with **12 records**

R Studio

db.rstudio.com / spark.rstudio.com

# Ideally, analyze in place, using **dplyr**

dplyr writes
the SQL
query

```sql
select count() from sales where
amount > 1000 group by month
```

Write dplyr:
```r
sales %>%
  filter(amount > 1000) %>%
  group_by(month) %>%
  tally()
```
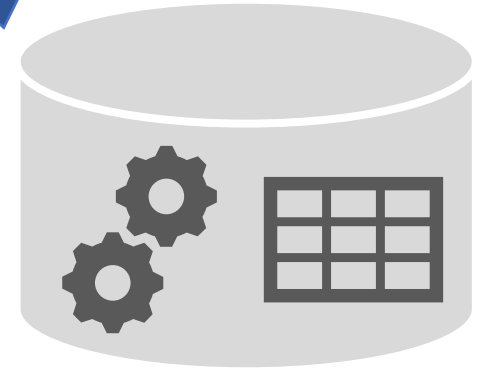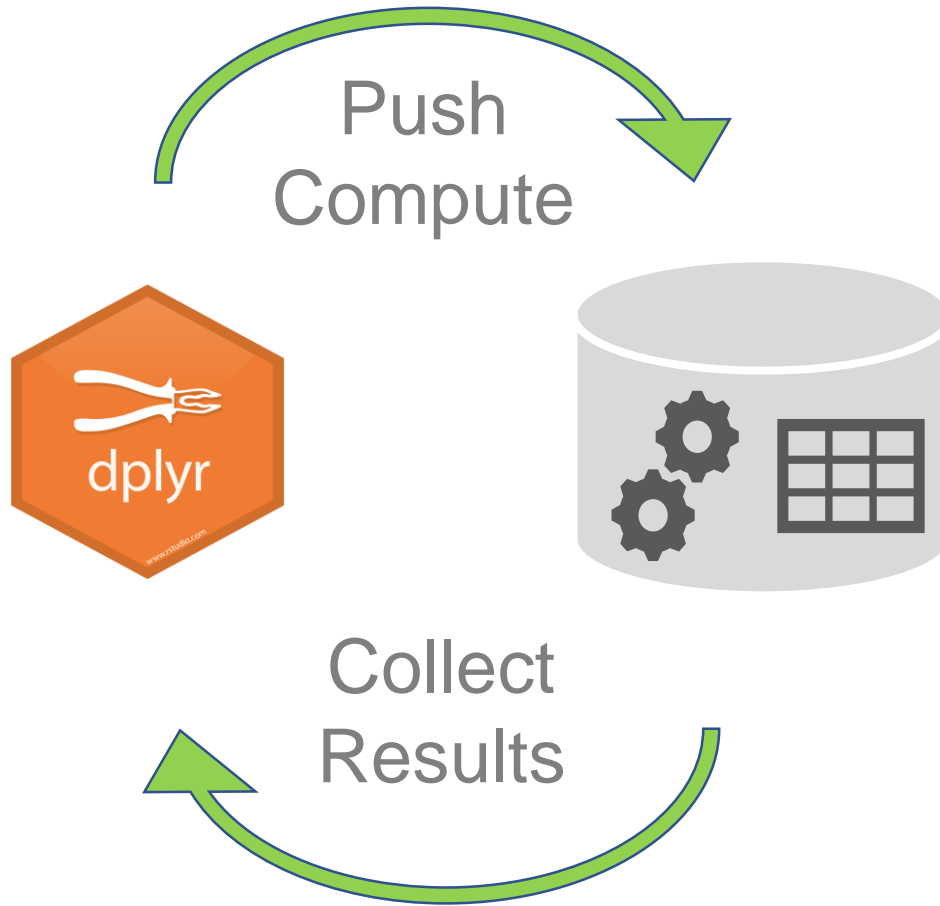
*"Number of sales over $1K by month"*

Returns a **data.frame** with **12 records**

R Studio

db.rstudio.com / spark.rstudio.com

# Available translations

1. Microsoft SQL Server
2. Oracle
3. Apache Hive
4. Apache Impala
5. PostgreSQL
6. MS Access
7. MariaDB (MySQL)
8. SQLite
9. Amazon Redshift
10. Teradata

# Advantages of using **dplyr**



Push Compute

Collect Results

1. dplyr translates to SQL

2. Take advantage of piped code

3. All your code is in R!

# Visualizations

R Studio

# Aggregate in DB, plot locally

**Single function**

**Local data**

**Aggregate**

**Plotting**

**Remote data**

R Studio

# dbplot for Histograms & Raster plots

**Single function**



dbplot          SQL                                    Bins

# Modeling with
**Databases**

Studio

# Option 1 - Modeling with a Database

**Process sample**

**Model & Test**

R Studio

# Score inside the DB using **tidypredict**

**Parse model**　　　　　**Score in database**

**SQL**

# Option 2 - Modeling in DB using **modeldb**

**Single Function**



modeldb      SQL               Model

# Modeling with **sparklyr**



Photo by Matthew Ronder-Seid on Unsplash

Studio

# Spark models (ML) available via **sparklyr**

| | |
|---|---|
| ALS | ml_als<br>ml_recommend<br>ml_als_factorization |
| Decision Trees | ml_decision_tree_classifier<br>ml_decision_tree<br>ml_decision_tree_regressor |
| Generalized Linear Regression | ml_generalized_linear_regression |
| Gradient Boosted Trees | ml_gbt_classifier<br>ml_gradient_boosted_trees<br>ml_gbt_regressor |
| K-Means Clustering | ml_kmeans<br>ml_compute_cost |
| Latent Dirichlet Allocation | ml_lda<br>ml_describe_topics<br>ml_log_likelihood<br>ml_log_perplexity<br>ml_topics_matrix |

| | |
|---|---|
| Linear Regression | ml_linear_regression |
| Logistic Regression | ml_logistic_regression |
| Multilayer Perceptron | ml_multilayer_perceptron_classifier<br>ml_multilayer_perceptron |
| Naive-Bayes | ml_naive_bayes |
| One Vs. Rest | ml_one_vs_rest |
| PCA (Estimator) | ft_pca ml_pca |
| Random Forest | ml_random_forest_classifier<br>ml_random_forest<br>ml_random_forest_regressor |
| Survival Regression | ml_aft_survival_regression<br>ml_survival_regression |

R Studio

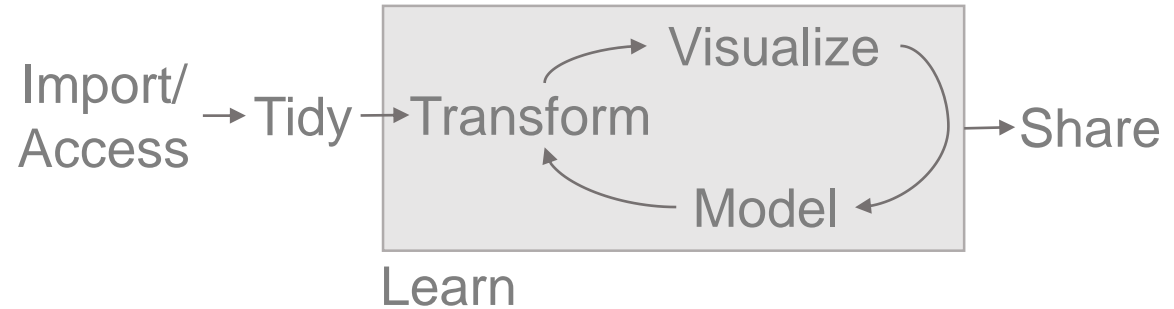# Production Pipelines



Photo by Iker Urteaga on Unsplash

Studio

# Different projects, different deliverables

## Data Science

- Deliverable: **Insights**
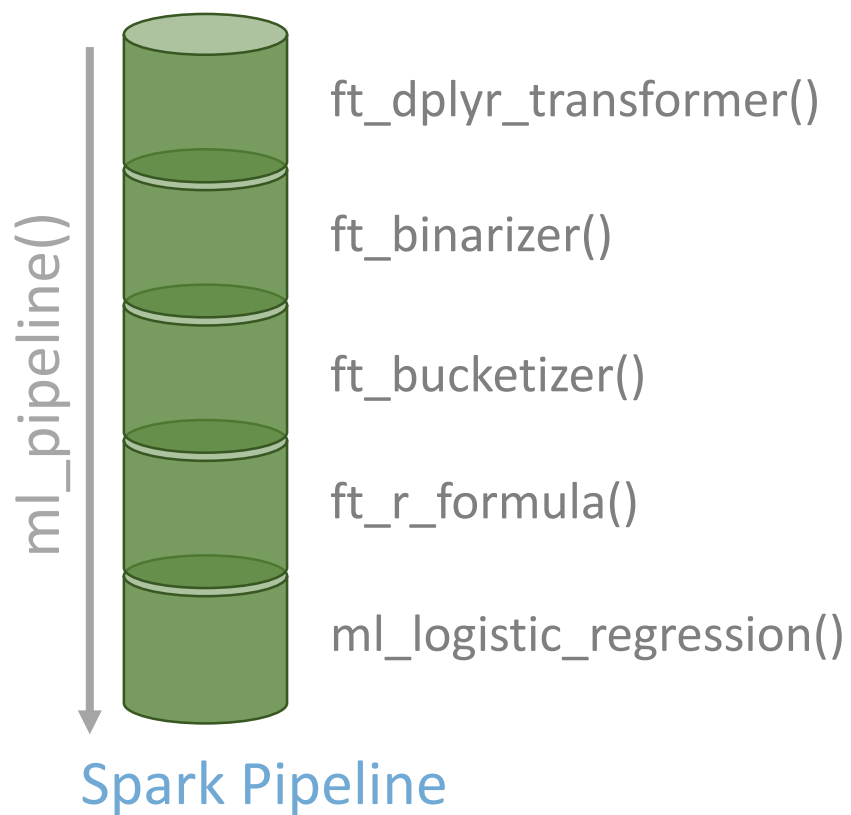- Experimental
- Iterative



Import/Access → Tidy → Transform → Visualize → Model → Share

Learn

## Production

- Deliverable: **Software**
- Tested
- Automated
- Apply SDLC

Report
*Ad-hoc | Emailed*

Predict
*Run daily*

Re-fit Model
*Run monthly*

# Spark pipelines types

**Estimator** (Plan)

**Transformer** (Fit)

ft_dplyr_transformer()

ft_binarizer()

ft_bucketizer()

ft_r_formula()

ml_logistic_regression()

ml_pipeline()

Spark Pipeline

Training dataset

ft_dplyr_transformer()

ft_binarizer()

ft_bucketizer()

ft_r_formula()

ml_logistic_regression()

ml_fit()

Spark PipelineModel

R Studio

db.rstudio.com / spark.rstudio.com

# Production Implementation



**DS Project**

**Frequency:** One time

**Frequency:** One time

**Fit Pipeline**

**Frequency:** Monthly

**PipelineModel**

**Frequency:** On-demand or daily

**Model Transform**

# Dashboards

R Studio

# Normal Shiny app



**Publish**

# Database + Shiny Dashboard

**Summary**

**Data Driven Dropdown**

10K

**Detail**

Reads data when dropdown changes

Reads data at app load

Reads data when user clicks on items

R Studio

db.rstudio.com / spark.rstudio.com

# R Tools for Big Data

| Access | Wrangle | Plot | Model | Automate |
|---|---|---|---|---|
| DBI<br>odbc<br>bigrquery<br>rpostgres<br>RMariaDB<br>monetdblite | dplyr (*via* dbplyr)<br>DBI<br>corrr (*in dev*) | ggplot2 (*via* dbplot)<br>corrr (*in dev*) | modeldb<br>tidypredict | tidypredict |
| sparklyr | sparklyr | | sparklyr<br>graphframes<br>rsparkling (H2O) | sparklyr<br>mleap |

Spark