# Compitiendo en Kaggle: Predicción de ventas

GRUPO DE R DE SEVILLA
Sevilla, 4 de Diciembre de 2018

**Javier Tejedor Aguilera**

✉ javier.tejedor.aguilera@gmail.com

in javier-tejedor-aguilera

**Datos**

# Rossmann Store Sales

Forecast sales using store, promotion, and competitor data

### Data Description

- **Store** - a unique Id for each store
- **Sales** - the turnover for any given day (this is what you are predicting)
- **Customers** - the number of customers on a given day
- **Open** - an indicator for whether the store was open: 0 = closed, 1 = open
- **StateHoliday** - indicates a state holiday. Normally all stores, with few exceptions, are closed on state holidays. Note that all schools are closed on public holidays and weekends. a = public holiday, b = Easter holiday, c = Christmas, 0 = None
- **SchoolHoliday** - indicates if the (Store, Date) was affected by the closure of public schools
- **StoreType** - differentiates between 4 different store models: a, b, c, d
- **Assortment** - describes an assortment level: a = basic, b = extra, c = extended
- **CompetitionDistance** - distance in meters to the nearest competitor store
- **CompetitionOpenSince[Month/Year]** - gives the approximate year and month of the time the nearest competitor was opened
- **Promo** - indicates whether a store is running a promo on that day
- **Promo2** - Promo2 is a continuing and consecutive promotion for some stores: 0 = store is not participating, 1 = store is participating
- **Promo2Since[Year/Week]** - describes the year and calendar week when the store started participating in Promo2
- **PromoInterval** - describes the consecutive intervals Promo2 is started, naming the months the promotion is started anew. E.g. "Feb,May,Aug,Nov" means each round starts in February, May, August, November of any given year for that store

| Store | StoreType | Assortment | CompetitionDistance | CompetitionOpenSinceMonth | CompetitionOpenSinceYear | Promo2 | Promo2SinceWeek | Promo2SinceYear | PromoInterval | DayOfWeek | Date | Sales | Open | Promo | StateHoliday | SchoolHoliday |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | c | a | 1270 | 9 | 2008 | 0 | NA | NA | | 2 | 01/01/2013 | 0 | 0 | 0 | a | 1 |
| 1 | c | a | 1270 | 9 | 2008 | 0 | NA | NA | | 3 | 02/01/2013 | 5530 | 1 | 0 | 0 | 1 |
| 1 | c | a | 1270 | 9 | 2008 | 0 | NA | NA | | 4 | 03/01/2013 | 4327 | 1 | 0 | 0 | 1 |
| 1 | c | a | 1270 | 9 | 2008 | 0 | NA | NA | | 5 | 04/01/2013 | 4486 | 1 | 0 | 0 | 1 |
| 1 | c | a | 1270 | 9 | 2008 | 0 | NA | NA | | 6 | 05/01/2013 | 4997 | 1 | 0 | 0 | 1 |
| | | | | | | | | | | | | | | | | |
| 1115 | d | c | 5350 | NA | NA | 1 | 22 | 2012 | Mar,Jun,Sept,Dec | 7 | 13/09/2015 | NA | 0 | 0 | 0 | 0 |
| 1115 | d | c | 5350 | NA | NA | 1 | 22 | 2012 | Mar,Jun,Sept,Dec | 1 | 14/09/2015 | NA | 1 | 1 | 0 | 0 |
| 1115 | d | c | 5350 | NA | NA | 1 | 22 | 2012 | Mar,Jun,Sept,Dec | 2 | 15/09/2015 | NA | 1 | 1 | 0 | 0 |
| 1115 | d | c | 5350 | NA | NA | 1 | 22 | 2012 | Mar,Jun,Sept,Dec | 3 | 16/09/2015 | NA | 1 | 1 | 0 | 0 |
| 1115 | d | c | 5350 | NA | NA | 1 | 22 | 2012 | Mar,Jun,Sept,Dec | 4 | 17/09/2015 | NA | 1 | 1 | 0 | 0 |

**Métrica de evaluación**

## Overview

**Description**

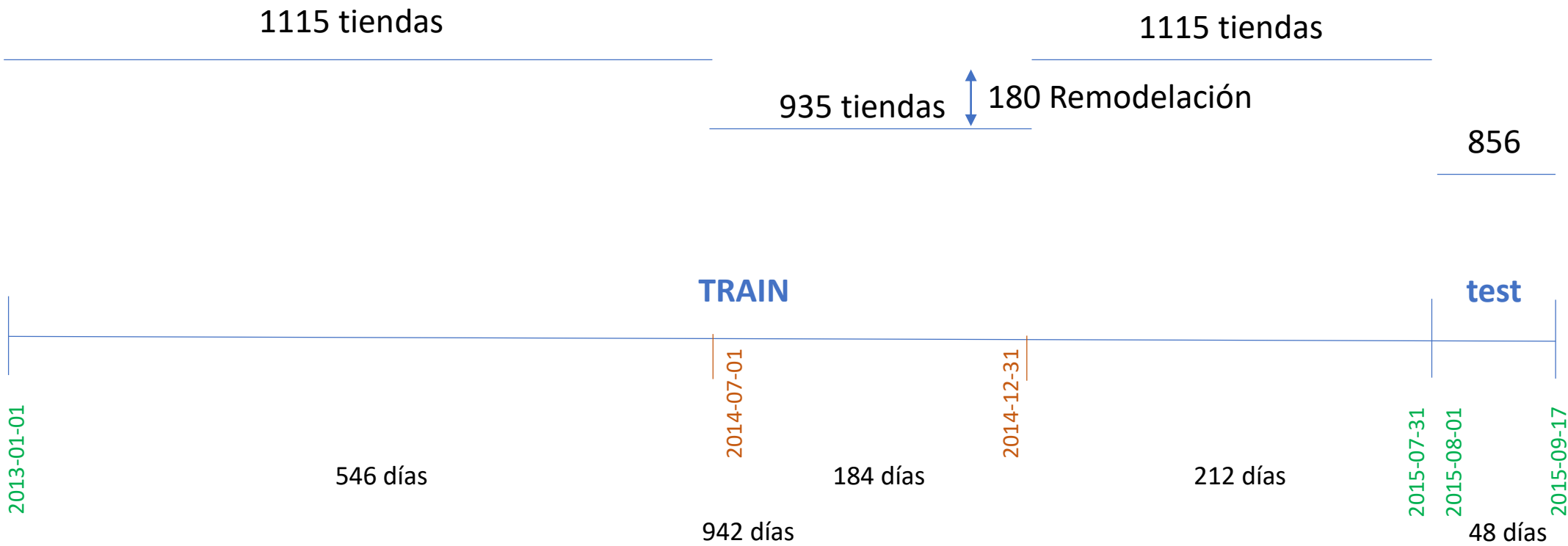**Evaluation**

**Prizes**

**Timeline**

Submissions are evaluated on the Root Mean Square Percentage Error (RMSPE). The RMSPE is calculated as

$$\text{RMSPE} = \sqrt{\frac{1}{n}\sum_{i=1}^{n}\left(\frac{y_i - \hat{y}_i}{y_i}\right)^2},$$

where y_i denotes the sales of a single store on a single day and yhat_i denotes the corresponding prediction. Any day and store with 0 sales is ignored in scoring.

Rossmann Store Sales

Forecast sales using store, promotion, and competitor data

1115 tiendas

1115 tiendas

935 tiendas ⇕ 180 Remodelación

856

TRAIN

test

2013-01-01

2014-07-01

2014-12-31

2015-07-31

2015-08-01

2015-09-17

546 días

184 días

212 días

942 días

48 días

**Estrategia de validación / crossvalidación**

**Rossmann Store Sales**

Forecast sales using store, promotion, and competitor data

TRAIN | Test

Crossvalidación 5-Fold Random

| Iteración 1 | Train | Train | Train | Train | Val | Test |
| Iteración 2 | Train | Train | Train | Val | Train | Test |
| Iteración 3 | Train | Train | Val | Train | Train | Test |
| Iteración 4 | Train | Val | Train | Train | Train | Test |
| Iteración 5 | Val | Train | Train | Train | Train | Test |

❌

Validación temporal

| Iteración 1 | Train | Val | Test |

✅

Crossvalidación 5-Fold temporal

| Iteración 1 | Train | Val | | | Test |
| Iteración 2 | Train | Val | | | Test |
| Iteración 3 | Train | Val | | Test |
| Iteración 4 | Train | Val | | Test |
| Iteración 5 | Train | Val | Test |

✅

tiempo

El dataset de validación debe mimetizar el test

**kaggle**

**ROSSMANN**

**Rossmann Store Sales**

Forecast sales using store, promotion, and competitor data

1115 tiendas

1115 tiendas

935 tiendas

856

**train**

**val**   **test**

2014-07-01

2014-12-31

2013-01-01

2015-06-13
2015-06-14
2015-07-31
2015-08-01
2015-09-17

546 días

184 días

164 días

894 días

48 días   48 días

ROSSMANN

**Ingeniería de Variables
(Modelo Gradient Boosting Decision Trees)**

kaggle

ROSSMANN

**Rossmann Store Sales**

Forecast sales using store, promotion, and competitor data

Tienda: T
Fecha: 16-03-2015
Variables: dia, sem, mes, anio

| | |
|---|---|
| dia | 16 |
| DayOfWeek | L |
| sem | 12 |
| mes | 3 |
| anio | 2015 |

L

16-03-2015

2013-07          2014-01          2014-07          2015-01

**Rossmann Store Sales**

Forecast sales using store, promotion, and competitor data

**Modelo**

*dmlc*
**XGBoost**

Level-wise tree growth

## GBDT Hyper Parameter Tuning

| Hyper Parameter | Tuning Approach | Range | Note |
|---|---|---|---|
| # of Trees | Fixed value | 100-1000 | Depending on datasize |
| Learning Rate | Fixed => Fine Tune | [2 - 10] / # of Trees | Depending on # trees |
| Row Sampling | Grid Search | [.5, .75, 1.0] | |
| Column Sampling | Grid Search | [.4, .6, .8, 1.0] | |
| Min Leaf Weight | Fixed => Fine Tune | 3/(% of rare events) | Rule of thumb |
| Max Tree Depth | Grid Search | [4, 6, 8, 10] | |
| Min Split Gain | Fixed | 0 | Keep it 0 |

Best GBDT implementation today: https://github.com/tqchen/xgboost
by **Tianqi Chen** (U of Washington)

DataRobot

| Xgboost | • Eta<br>• Gamma<br>• Max_depth<br>• Min_child_weight<br>• Subsample<br>• Colsample_bytree<br>• Lambda<br>• alpha | • 0.01,0.015, 0.025, 0.05, 0.1<br>• 0.05-0.1,0.3,0.5,0.7,0.9,1.0<br>• 3, 5, 7, 9, 12, 15, 17, 25<br>• 1, 3, 5, 7<br>• 0.6, 0.7, 0.8, 0.9, 1.0<br>• 0.6, 0.7, 0.8, 0.9, 1.0<br>• 0.01-0.1, 1.0 , RS*<br>• 0, 0.1, 0.5, 1.0 RS* |
|---|---|---|

**Ideas de mejora**

## 1. Generar N modelos, uno por día a predecir

Idea propuesta por **Danijel Kivaranovic** (Kaggle Grandmaster) en la competición
Recruit Restaurant Visitor Forecasting de Kaggle y que expuso en Kaggle Days Varsovia.

Recruit Restaurant Visitor Forecasting
Predecir los visitantes en 314 restaurantes en los próximas 39 días

Nota: Imagen extraída de la presentación de **Danijel Kivaranovic** en Kaggle Days Warsaw

39 días -> 39 Modelos -> 39 datasets de validación -> 39 scores de validación

Como era de esperar, cuanto más lejana en el tiempo la predicción, mayor es el error.

**Más ideas de mejora**
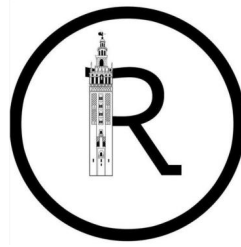
1. Mejorar ingeniería de variables
   - Mejorar algunas de las variables calculadas
   - Cambiar codificación de las categóricas: oneHotEncoding, LabelEncoding, Frecuencia, etc.
   - Nuevas variables

2. Bagging
   - Promediar modificando algunos hiperparámetros
   - Promediar modificando la forma de cálculo de algunas variables (cambiando codificación categóricas, etc.)

3. Stacking
   - Combinar las predicciones de varios modelos (RF, RNN, LR, ARIMA, …) sacando lo mejor de cada uno. Cuanto mayor diversidad, mejor.

**Muchas gracias
por vuestra atención!**