

Advanced Industrial Turnover Index using Statistical Learning Algorithms

The Use of R in Official Statistics uRos 2021

S. Barragán¹, L. Barreñada³, J.F. Calatrava¹,
J.C. Gálvez Sáenz de Cueto¹, J.M. Martín del Moral²,
E. Rosa-Perez^{2,3}, D. Salgado^{1,3}

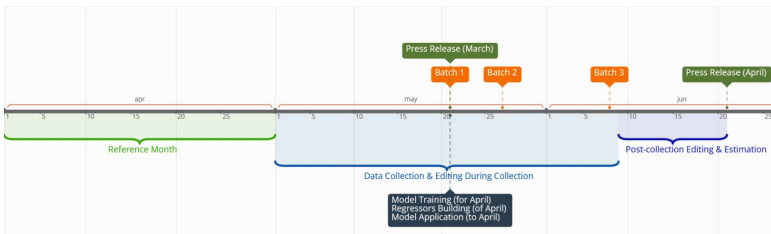
¹Dept. Methodology and Development of Statistical Production, Statistics Spain (INE)

²S.G. for Industrial and Services Statistics, Statistics Spain (INE)

³Dept. Statistics and Operations Research, Complutense University of Madrid

Bucharest, 24-26 November 2021

- (i) Short-Term Business Statistics: monthly release.
- (ii) Accuracy vs. Timeliness



- Monthly fixed-base Laspeyres indices with cut-off sampling of establishments for total industrial turnover for each dissemination cell U_d

$$Y_{U_d}^{(m)} = \sum_{k \in U_d} y_k^{(m, \text{val})} \rightsquigarrow I_d^{(m)}, \Delta_d^{(m, m-12)}, \Delta_d^{(m, m-1)}$$

- Tested advanced version:
 r_t subsample collected up to time $t < t_{\text{release}}(m)$

$$Y_{U_d}^{(m)}(t) = \sum_{k \in r_{t,d}} y_{kt}^{(m, \text{ed})} + \sum_{k \in U_d - r_{t,d}} \hat{y}_{kt}^{(m, \text{val})}$$

Prediction method: Statistical Learning

- Used data:
 - Survey data from ITI.
 - Survey data from IPI and IPRI.
- Target Variable: validated total turnover $y_k^{(m, val)}$ per unit and reference month.

- Regressors (287):

	ID	Cross	Long	Cross+ Long	External
Hist. Series	✓	✓	✓	✗	✗
Running Month	✓	✓	✗	✓	✓

- Light Gradient Boosting on Regression Trees.

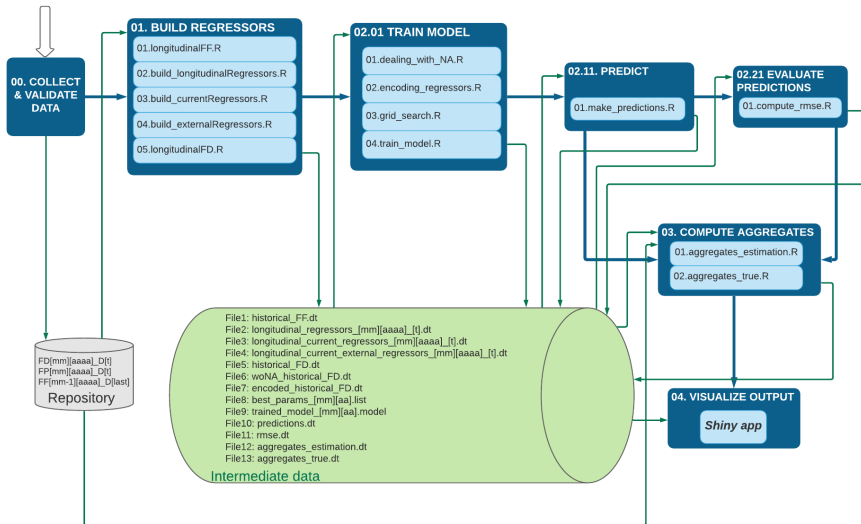
Regressors Importance

Overview


Methodology

Results



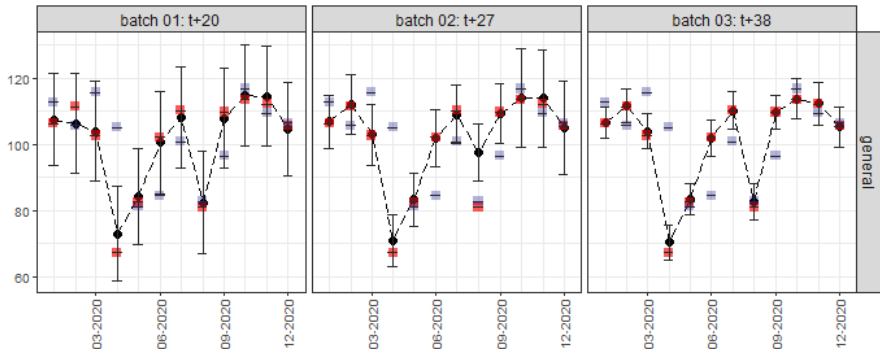


Pilot: Software and execution

-  Packages Stack:
 - **Read and data management**
data.table, stringr, stringi, haven,
lubridate, xml2, fastReadfwf⁽¹⁾, RepoTime⁽¹⁾,
StQ⁽¹⁾, RepoReadWrite⁽¹⁾, RepoUtils⁽¹⁾.
 - **Statistical Learning**
mltools, lightgbm, xgboost, psych.
 - **Visualization**
shinydashboard, shiny, shinythemes, ggthemes,
gridExtra, latex2exp, ggplot2.
 - Windows10, 32GB RAM, Intel(R) Core™ i7-4790 CPU 3.60GHz.
 - Executed on more than 60 consecutive months (12000 establishments per month).
- (1): In-house developed.

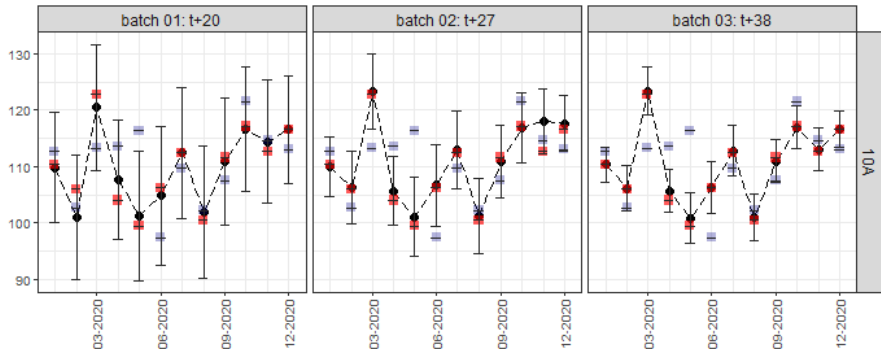
Industrial Turnover Index

Index Version ● advanced ± rmse ■ final ▣ initial



Industrial Turnover Index

Index Version ● advanced ± rmse ■ final ▣ initial



Overview
Method
Result:

- Predict semicontinuous variables \rightsquigarrow seasonality
- Predict measurement errors: $y_{kt}^{(m,ed)} - y_k^{(m,val)}$
- Model residuals \rightsquigarrow uncertainty intervals
- Implement in production \rightsquigarrow interoperable software tools
- Generalize to other STS (probabilistic sampling).
- Daily Data Collection and daily turnover availability in the reference month.