

EST-25134: Aprendizaje Estadístico

Profesor: Alfredo Garbuno Iñigo — Primavera, 2023 — Interpretabilidad y explicabilidad de modelos predictivos.

Objetivo: En aplicaciones de modelos predictivos usualmente se consideran modelos con alto poder predictivo. A su vez, estos modelos son altamente complejos y es difícil *explicar* el cómo una predicción es realizada a consecuencia del vector de atributos en consideración. En esta sección estudiaremos algunas de las nociones de interpretabilidad de modelos.

Lectura recomendada: Los libros de Biecek and Burzykowski [1] y Molnar [2] son dos publicaciones recientes que tratan temas de interpretabilidad y explicabilidad de modelos.

1. INTRODUCCIÓN

En aplicaciones de modelos predictivos usualmente se consideran modelos con alto poder predictivo. A su vez, estos modelos son altamente complejos y es difícil *explicar* el cómo una predicción es realizada a consecuencia del vector de atributos en consideración. En esta sección estudiaremos algunas de las nociones de interpretabilidad de modelos.

Para algunos modelos, como regresión lineal o árboles de decisión, es relativamente sencillo interpretar las relaciones entre atributos y variable respuesta.

Para ilustrar retomaremos el ejemplo de productos de Ikea, el cual es original de: [Tune random forests for #TidyTuesday IKEA prices](#).

Los datos que tenemos disponibles son los siguientes.

```
1 ikea_df <- ikea >
2   select(price, name, category, depth, height, width) >
3   mutate(price = log10(price)) >
4   mutate_if(is.character, factor)
5
6 ikea_df > print(n = 5)
```

```
1 # A tibble: 3,694 × 6
2   price name                category    depth height width
3   <dbl> <fct>                <fct>    <dbl>  <dbl> <dbl>
4 1  2.42 FREKVEN          Bar furniture    NA     99    51
5 2  3.00 NORDVIKEN        Bar furniture    NA    105    80
6 3  3.32 NORDVIKEN / NORDVIKEN Bar furniture    NA     NA    NA
7 4  1.84 STIG              Bar furniture    50    100    60
8 5  2.35 NORBERG          Bar furniture    60     43    74
9 # ... with 3,689 more rows
10 # Use 'print(n = ...)' to see more rows
```

Los cuales son sometidos a nuestro típico flujo de trabajo de ajuste de modelos predictivos junto con un proceso de separación de muestras para métricas de generalización y selección de hiper-parámetros.

```

1 set.seed(123)
2 ikea_split <- initial_split(ikea_df, strata = price)
3 ikea_train <- training(ikea_split)
4 ikea_test <- testing(ikea_split)
5
6 set.seed(234)
7 ikea_folds <- vfold_cv(ikea_train, strata = price)

```

```

1 library(textrecipes)
2 ranger_recipe <-
3   recipe(formula = price ~ ., data = ikea_train) ▷
4   step_other(name, category, threshold = 0.01) ▷
5   step_clean_levels(name, category) ▷
6   step_impute_knn(depth, height, width)

```

```

1 linear_recipe <-
2   recipe(formula = price ~ ., data = ikea_train) ▷
3   step_other(name, category, threshold = 0.01) ▷
4   step_clean_levels(name, category) ▷
5   step_impute_knn(depth, height, width) ▷
6   step_dummy(all_nominal_predictors()) ▷
7   step_normalize(all_predictors())

```

1.1. Especificación del modelo

```

1 linear_spec <-
2   linear_reg(penalty = 1e-3) ▷
3   set_mode("regression") ▷
4   set_engine("glmnet")
5
6 linear_workflow <-
7   workflow() ▷
8   add_recipe(linear_recipe) ▷
9   add_model(linear_spec)

```

```

1 ranger_spec <-
2   rand_forest(trees = 1000) ▷
3   set_mode("regression") ▷
4   set_engine("ranger")
5
6 ranger_workflow <-
7   workflow() ▷
8   add_recipe(ranger_recipe) ▷
9   add_model(ranger_spec)

```

```

1 all_cores <- parallel::detectCores(logical = TRUE) - 1
2 library(doParallel)
3 cl <- makePSOCKcluster(all_cores)
4 registerDoParallel(cl)

```

```

1 ikea_lm <- linear_workflow ▷ fit(data = ikea_train)
2 ikea_rf <- ranger_workflow ▷ fit(data = ikea_train)

```

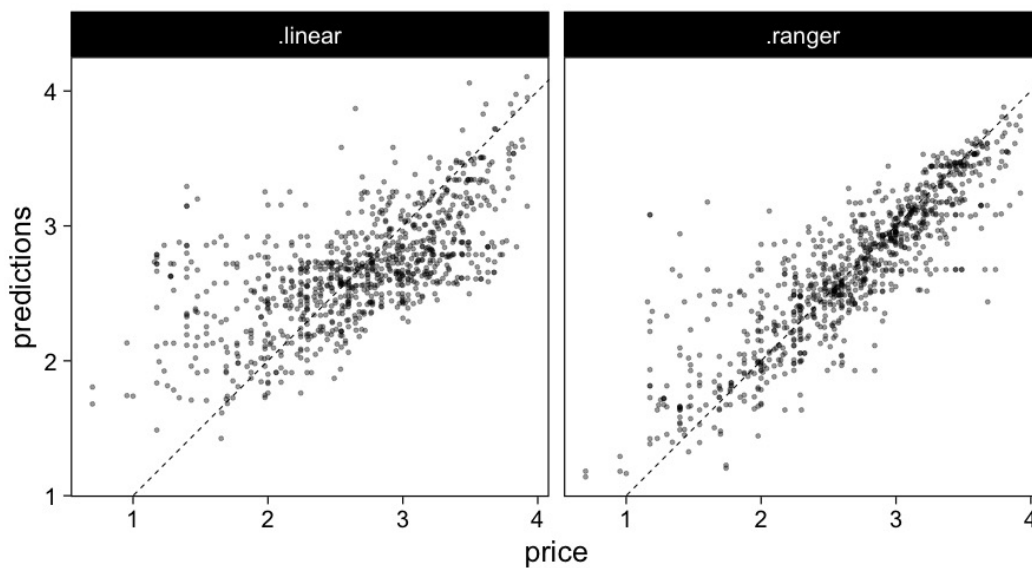
2. INTERPRETABILIDAD

Iremos explorando los conceptos necesarios para interpretabilidad conforme los necesitemos. Primero necesitaremos herramientas de trabajo desde R, y para esta tarea podemos usar `lime`, `vip` y `DALEXtra`.

En general podemos usar:

- `vip` para usar métodos basados en algún modelo en particular para aprovechar la estructura del modelo predictivo.
- `DALEX` para usar métodos que no requieren de una estructura en particular (usaremos `DALEXtra` para compatibilidad con `tidymodels`).

```
1 library(DALEXtra)
```



Para poder comenzar lo que tenemos que hacer es crear los objetos de `DALEX` (`moDel Agnostic Language for Exploration and eXplanation`).

```
1 explainer_lm <-  
2   explain_tidymodels(  
3     ikea_lm,  
4     data = ikea_train > select(-price),  
5     y     = ikea_train > pull(price),  
6     label = "linear model",  
7     verbose = FALSE  
8   )
```

```
1 explainer_rf <-  
2   explain_tidymodels(  
3     ikea_rf,  
4     data = ikea_train > select(-price),  
5     y     = ikea_train > pull(price),  
6     label = "random forest",  
7     verbose = FALSE  
8   )
```

3. MÉTODOS DE INTERPRETABILIDAD LOCAL

Los siguientes métodos que veremos son `métodos locales` es decir, tomamos una $x_0 \in \mathcal{X} \subset \mathbb{R}^p$ en particular y exploramos la respuesta a partir de este punto. Por ejemplo, consideremos como x_0 la localidad donde queremos explorar el modelo.

```
1 set.seed(123)
2 mueble <- ikea_test > sample_n(1)
3 mueble
```

Sabemos de modelos lineales que los coeficientes están asociados a las contribuciones de cada predictor a la respuesta. Usualmente, interpretados bajo un principio *ceteris paribus* (interpretado en nuestro contexto: dejando constantes los demás predictores).

```
1 ikea_lm > extract_fit_parsnip() >
2   tidy() >
3   print(n = 5)
```

3.0.1. Para pensar: Un profesional de la estadística les recordaría el concepto de *ceteris paribus* en el contexto de regresión. Es alrededor del vector $0 \in \mathcal{X}$ el que usamos para la interpretación o es alrededor del individuo promedio $\bar{x} \in \mathcal{X}$ el que usamos para interpretar el ajuste?

4. EXPANSIONES LINEALES LOCALES

Una vez que hemos decidido sobre cual individuo (observación o instancia) queremos hacer la expansión podemos usar **DALEX** para poder crear métricas de sensibilidad de cambios del valor promedio de la predicción derivado de cambios individuales en los atributos.

```
1 lm_breakdown <- predict_parts(explainer = explainer_lm, new_observation =
  mueble)
2 lm_breakdown
```

Lo mismo podemos hacer para nuestro modelo de **random forest**. En este tipo de tablas interpretamos cómo cada cambio va alejandonos de nuestro *intercepto* (la respuesta promedio de nuestro modelo predictivo).

```
1 rf_breakdown <- predict_parts(explainer = explainer_rf, new_observation =
  mueble)
2 rf_breakdown
```

La interpretación cambia de acuerdo al orden en como se van presentando los cambios en los atributos y para esto podemos usar el modelo lineal como una heurística de orden.

```
1 predict_parts(explainer = explainer_rf,
2               new_observation = mueble,
3               order = lm_breakdown$variable_name)
```

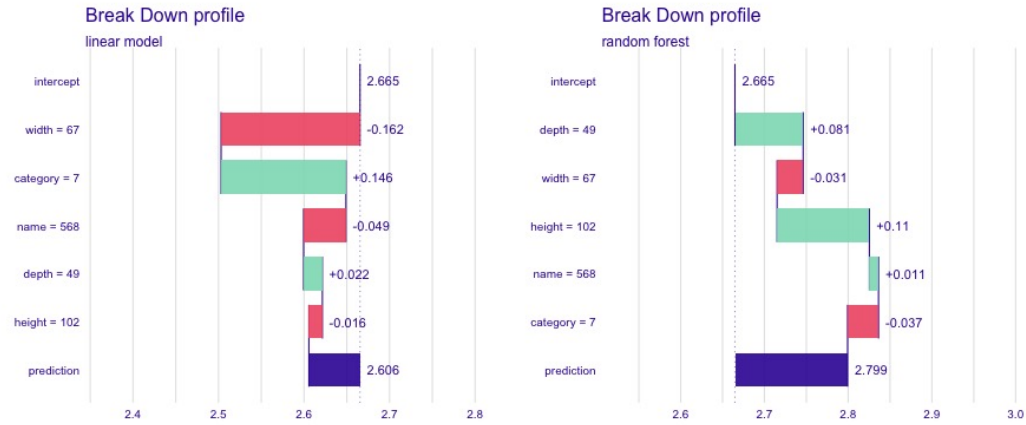
```
1               contribution
2 random forest: intercept      2.665
3 random forest: width = 67    -0.062
4 random forest: category = 7  -0.049
```

5 SHAP VALUES

```

5 random forest: name = 568          -0.027
6 random forest: depth = 49          0.183
7 random forest: height = 102        0.090
8 random forest: prediction          2.799

```

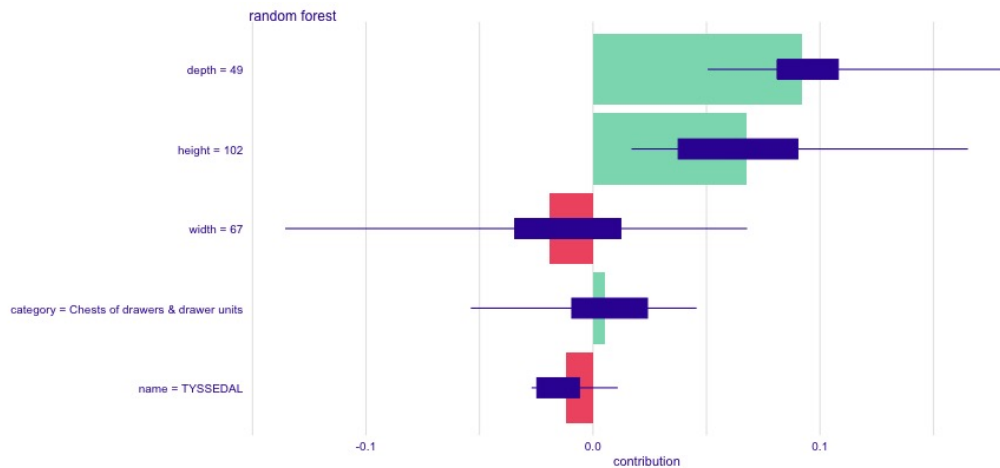


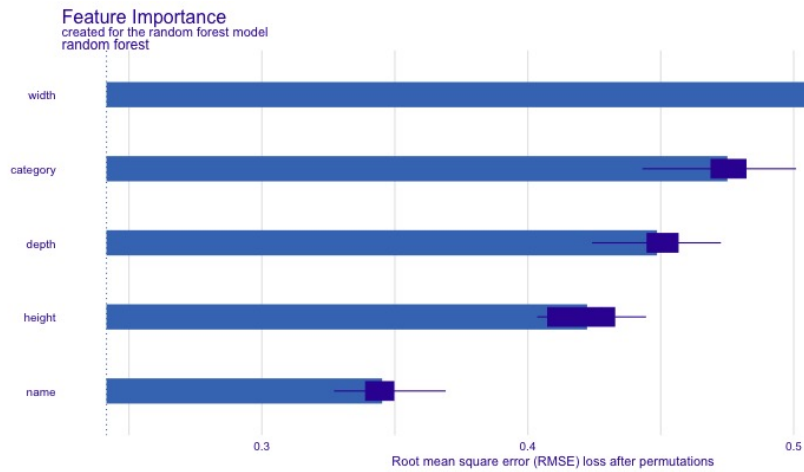
5. SHAP VALUES

```

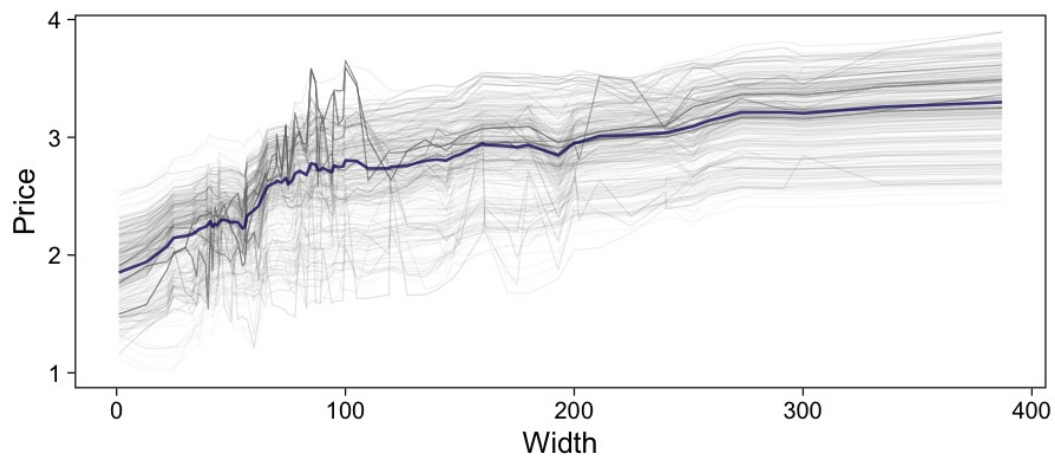
1 set.seed(1801)
2 shap_mueble <-
3   predict_parts(
4     explainer = explainer_rf,
5     new_observation = mueble,
6     type = "shap",
7     B = 20
8   )

```

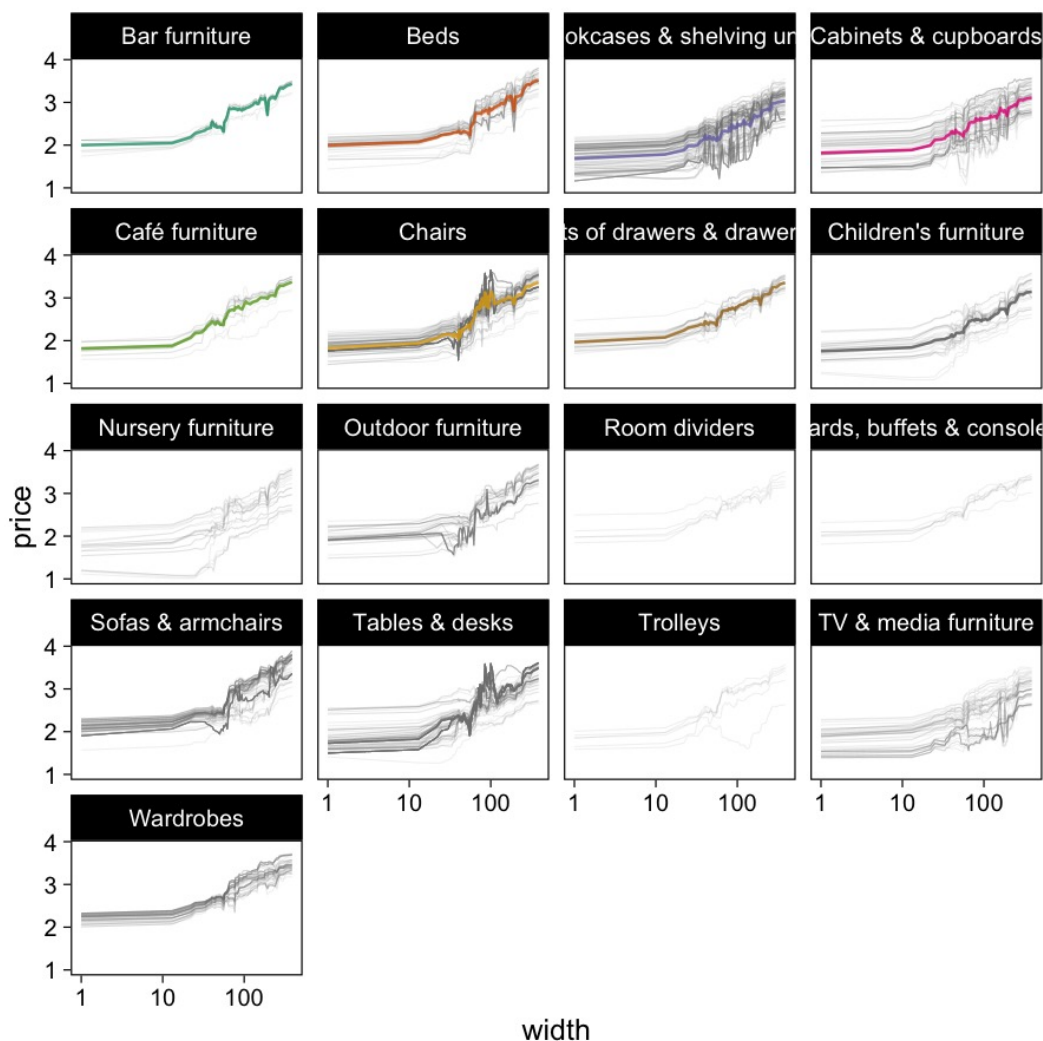




```
1 set.seed(1805)
2 pdp_width <- model_profile(explainer_rf, N = 500, variables = "width")
```



```
1 set.seed(1806)
2 pdp_wcat <- model_profile(explainer_rf, N = 1000,
3                           variables = "width",
4                           groups = "category")
```



REFERENCIAS

- [1] P. Biecek and T. Burzykowski. *Explanatory Model Analysis: Explore, Explain, and Examine Predictive Models*. Chapman & Hall/CRC Data Science Series. CRC Press, Boca Raton, first edition, 2021. ISBN 978-0-367-13559-1. [1](#)
- [2] C. Molnar. *Interpretable Machine Learning*. Lean Pub, 2020. [1](#)