

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", importance = "impurity")
5
6 ranger_workflow <-
7   workflow() ▷
8   add_recipe(ranger_recipe) ▷
9   add_model(ranger_spec)

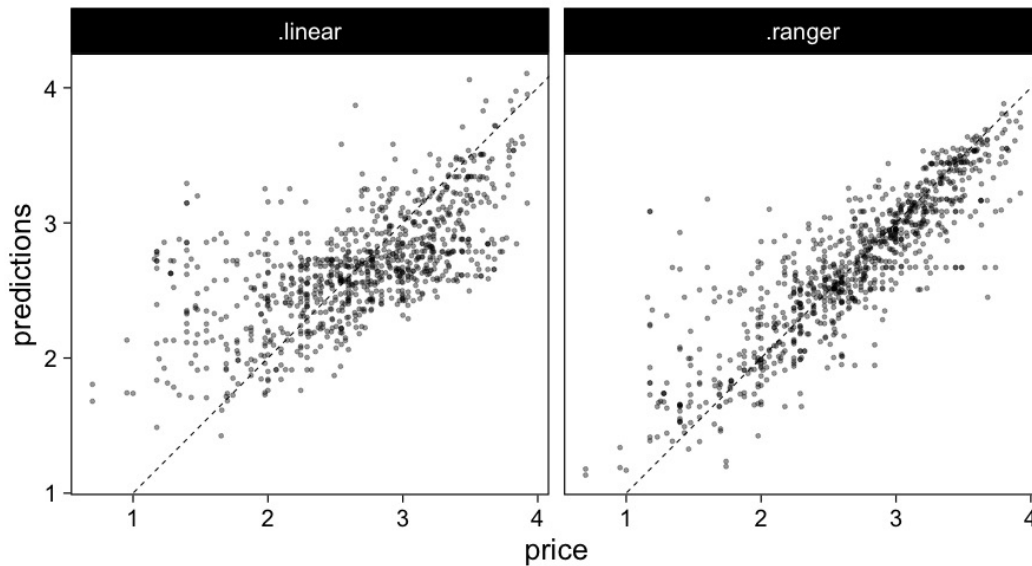
```

```

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 $x^* \in \mathcal{X} \subset \mathbb{R}^p$ en particular y exploramos la respuesta a partir de este punto. Por ejemplo, consideremos como x^* la observación donde queremos explorar el modelo.

```

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

```

```

1 # A tibble: 1 × 6
2   price name      category      depth height width
3   <dbl> <fct>    <fct>      <dbl>  <dbl> <dbl>
4 1   2.98 TYSSDAL Chests of drawers & drawer units    49    102    67

```

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 constantes).

```

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

```

```

1 # A tibble: 35 × 3
2   term      estimate penalty
3   <chr>      <dbl>    <dbl>
4 1 (Intercept)  2.67      0.001
5 2 depth        0.104      0.001
6 3 height        0.155      0.001
7 4 width         0.237      0.001
8 5 name_bekant  0.00497     0.001
9 # ... with 30 more rows
10 # Use 'print(n = ...)' to see more rows

```

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 $x^* \in \mathcal{X}$ 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. DESCOMPOSICIÓN SECUENCIAL ADITIVA

Recordemos que nuestras predicciones (en regresión) se pueden asociar a la función de regresión

$$\hat{f}(x) = \mathbb{E}[y|x], \quad (1)$$

siempre y cuando utilicemos pérdida cuadrática para realizar el ajuste.

Por ejemplo en regresión lineal podemos calcular el valor esperado de la respuesta para una observación x^* por medio de

$$\mathbb{E}[y|x^*] = \beta_0 + \beta_1 x_1^* + \cdots + \beta_p x_p^*, \quad (2)$$

donde los coeficientes β se ajustan por MCO.

En general, podemos calcular cómo cambia el valor esperado de y condicionado en que el atributo j tiene un valor de X_j por medio de

$$\iota(j, x^*) = \mathbb{E}[Y|x^*] - \mathbb{E}_{X_j}[\mathbb{E}[y|X_j]], \quad (3)$$

donde $\iota(j, x^*)$ mide la importancia de la variable j evaluada en el punto x^* .

Por ejemplo, en regresión lineal tenemos la expresión particular de

$$\iota(j, x^*) = \beta_j (x_j^* - \mathbb{E}[X_j]), \quad (4)$$

que podemos utilizar para expresar

$$\hat{f}(x^*) = (\text{prediccion media}) + \sum_{j=1}^p \iota(j, x^*). \quad (5)$$

En general, cuando usamos modelos no lineales podemos pensar en

$$\iota(j, x^*) = \mathbb{E}[Y|x_{1:j}^*] - \mathbb{E}[y|x_{1:j-1}^*], \quad (6)$$

para preservar una suma telescópica como la anterior.

Nota como el orden de los atributos afecta la descomposición de la predicción en términos individuales. Como es de esperarse es un mal resumen cuando hay interacción entre atributos.

Una vez que hemos decidido sobre cual individuo (observación o instancia) queremos hacer la expansión podemos usar **DALEX** para poder crear métricas de sensibilidad. Para esto utilizamos la función `predict_parts()`.

```
1 lm_breakdown <- predict_parts(  
2   explainer = explainer_lm,  
3   new_observation = mueble  
4 )  
5 lm_breakdown
```

```
1               contribution  
2 linear model: intercept      2.665  
3 linear model: width = 67    -0.162  
4 linear model: category = 7   0.146  
5 linear model: name = 568     -0.049  
6 linear model: depth = 49     0.022  
7 linear model: height = 102   -0.016  
8 linear model: prediction     2.606
```

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

```

1 rf_breakdown <- predict_parts(
2   explainer = explainer_rf,
3   new_observation = mueble
4 )
5 rf_breakdown

```

```

1               contribution
2 random forest: intercept      2.665
3 random forest: depth = 49     0.082
4 random forest: width = 67    -0.037
5 random forest: height = 102   0.111
6 random forest: name = 568     0.008
7 random forest: category = 7  -0.033
8 random forest: prediction     2.795

```

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.

```

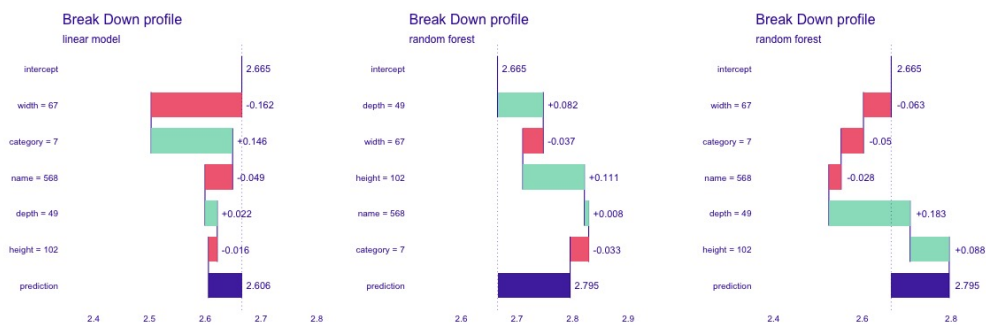
1 rfor_breakdown <- predict_parts(
2   explainer = explainer_rf,
3   new_observation = mueble,
4   order = lm_breakdown$variable_name
5 )
6 rfor_breakdown

```

```

1               contribution
2 random forest: intercept      2.665
3 random forest: width = 67    -0.063
4 random forest: category = 7  -0.050
5 random forest: name = 568    -0.028
6 random forest: depth = 49     0.183
7 random forest: height = 102   0.088
8 random forest: prediction     2.795

```



Podemos utilizar también la siguiente opción para explorar posibles contribuciones derivadas de interacciones.

```

1 rfin_breakdown <- predict_parts(
2   explainer = explainer_rf,
3   new_observation = mueble,
4   type = "break_down_interactions"

```

```
5 )
6 rfin_breakdown
```

```
1
2
3
4
5
6
7
      contribution
random forest: intercept      2.665
random forest: depth = 49      0.082
random forest: width:category = 67:7 -0.033
random forest: height = 102     0.090
random forest: name = 568      -0.009
random forest: prediction      2.795
```

5. EXPANSIONES LINEALES LOCALES

Podemos explorar la idea de aproximar un modelo predictivo sumamente complejo por uno altamente transparente. Esta es la idea detras de LIME (*Local Interpretable Model-agnostic Explanations*).

Por ejemplo, podemos utilizar la función `predict_surrogate()` de la siguiente manera

```
1 xgb_spec <-
2   boost_tree(trees = 1000) >
3   set_mode("regression") >
4   set_engine("xgboost")
5
6 xgb_workflow <-
7   workflow() >
8   add_recipe(ranger_recipe > step_dummy(all_nominal_predictors())) >
9   add_model(xgb_spec)
```

```
1 ikea_xgb <- xgb_workflow > fit(data = ikea_train)
```

```
1 explainer_xgb <-
2   DALEX::explain(
3     ikea_xgb,
4     data = ikea_train > select(-price),
5     y     = ikea_train > pull(price),
6     label = "boosted trees",
7     verbose = FALSE
8   )
```

```
1 library(lime)
2 set.seed(108)
3 model_type.dalex_explainer <- DALEXtra::model_type.dalex_explainer
4 predict_model.dalex_explainer <- DALEXtra::predict_model.dalex_explainer
5
6 lime_mueble <- predict_surrogate(
7   explainer = explainer_xgb,
8   new_observation = mueble,
9   n_features = 3,
10  n_permutations = 500,
11  type = "lime"
12 )
13
14 lime_mueble > print(width = 85)
```

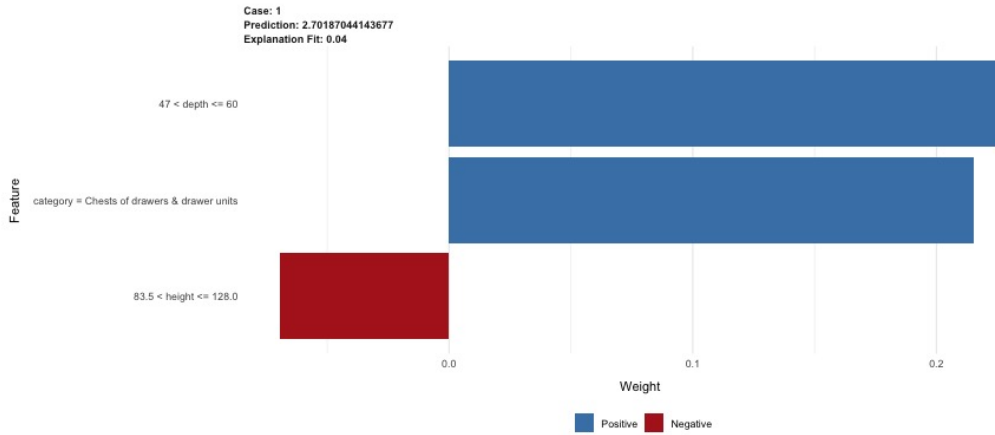
```

1 # A tibble: 3 × 11
2   model_type case model_r2 model_intercept model_prediction feature
3   <chr>      <chr>    <dbl>      <dbl>          <dbl> <chr>
4 1 regression 1      0.0401      2.40          2.77 depth
5 2 regression 1      0.0401      2.40          2.77 category
6 3 regression 1      0.0401      2.40          2.77 height
7   feature_value feature_weight feature_desc      data ...1
8     <dbl>          <dbl> <chr>          <list>      <dbl>
9 1         49      0.226 47 < depth ≤ 60    <named list>  2.70
10 2          7      0.215 category = Chests of ...dra <named list>
11 3        102     -0.0694 83.5 < height ≤ 128.0    <named list>  2.70
12 # ... with abbreviated variable name 1prediction

```

El modelo lineal descrito arriba realiza predicciones por medio de

$$\hat{f}(x_{\text{mueble}}) = 2.40 + 0.226 \times x_1 + 0.215 \times x_2 - 0.07 \times x_3. \quad (7)$$



5.1. Construcción de LIME

La idea es sencilla. Tenemos un modelo predictivo $\hat{f}(x)$ que hemos ajustado con un conjunto de datos. Ahora lo que buscamos es un modelo *transparente* g tal que

$$\hat{g}(x) = \arg \min_{g \in \mathcal{G}^*} \|g - \hat{f}\|^2 + \Omega(g), \quad (8)$$

donde $\mathcal{G}^* = \{g : \mathcal{N}(x^*) \subset \mathcal{X} \rightarrow \mathcal{Y} | g \text{ es una función sencilla} \}$, $\|\cdot\|$ es una norma apropiada y $\Omega(\cdot)$ asigna una penalización por complejidad de g .

- Usualmente restringimos \mathcal{G}^* a ser un espacio de funciones lineales que incluye sólo d entradas (con $d \ll p$).
- Para generar puntos en $\mathcal{N}(x^*)$ creamos perturbaciones a partir de x^* y evaluamos el modelo entrenado \hat{f} en esos puntos.
- Ajustamos un modelo de regresión lineal para el conjunto de datos sintéticos.

5.2. Observaciones

- No hacemos supuestos sobre el modelo \hat{f} .
- La representación en menores dimensiones nos ayuda a mantener los atributos en un marco manejable.

- La aproximación sólo es local.
- Puede y ha sido utilizado en modelos de texto y visión por computadora.
- Hay que recordar que sólo es una interpretación del modelo ajustado, no de los datos.
- Nosotros utilizamos sólo las funciones de `lime` pero también pueden utilizar `iml` o `localModels`, pueden ver mas [aquí](#).

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](#)