

---

UNIVERSIDAD DE BUENOS AIRES  
FACULTAD DE INGENIERÍA  
Año 2023 - 1<sup>er</sup> Cuatrimestre

ORGANIZACIÓN DE DATOS (75.06 /95.58)

TRABAJO PRACTICO 1: RESERVAS DE HOTEL

CHECKPOINT 2: Árboles de decisión  
FECHA: 28/04/23

INTEGRANTES:

Gonzalez, Tomás	#108193
<togonzalez@fi.uba.ar>	
Moreno del Ruvo, Valentina	#107948
<vmr18@gmail.com>	
Pol, Juan Manuel	#108448
<jpol@fi.uba.ar>	

### Introducción

El objetivo de este segmento fue generar predicciones de la variable target *is\_canceled* a partir de la creación y análisis de árboles de decisiones, y la optimización de sus hiperparámetros.

### Transformación de variables categóricas

Decodificamos variables **dummies** a partir de cada valor de variables categóricas como *agent* o *meal* de manera que los árboles puedan trabajar con estas. Para el caso de la variable *company*, decidimos pasarla a la binaria *no\_company*, que refiere a si hay o no una compañía responsable de la reservación.

### Elección de variables

Dada la cantidad de valores únicos de las variables *agent* y *country*, decidimos generar 2 árboles (uno de baja profundidad y otro de mayor pero podado) para 3 casos distintos: **Sin ninguna de las dos**, **Solo con *country***, **Solo con *agent***. Y elegir la variable con la que mejor predicción logra el árbol. Una vez probados todos los casos anteriormente descriptos, concluimos que el mejor caso era el dataset **solo con *country*(podado)** (priorizando el *f-score*). Además, notamos que el valor 'portugal' era el único que incidía significativamente en el árbol, por lo que optamos por dejar una sola variable *is\_portugal* que refiera a si provienen de portugal o no.

### Obtención del modelo

Una vez elegido el conjunto de variables para entrenar nuestro modelo, pasamos a la búsqueda de hiperparámetros, utilizando **Random search** para obtener las combinaciones posibles al azar, y luego **10-fold CV** para medir dichas combinaciones. Nos quedamos con la combinación con mayor desempeño.

### Testeo del modelo

Una vez que tenemos el mejor modelo, lo entrenamos con los datos de *train* para luego 'testearlo' con los datos de *test*. Obtenemos los resultados y subimos el *score*, de 0.81017, a la competencia de Kaggle.