

Maximizando el ROI de los datos mediante Imputación de valores perdidos con *Ensamble Learning*

Germán Rosati
german.rosati@gmail.com

UNTREF / MTEySS / Digital House

26 de Abril de 2017

- 1 ¿Qué es y como se genera un dato perdido?
- 2 ¿Cómo lidiar con los datos perdidos?
 - Técnicas tradicionales (imputación simple)
- 3 Ejercicio de aplicación

¿Qué es un dato perdido?

¿Qué es un valor perdido?

- Los datos son caros
 - Encuestas
 - SRMs
 - Logs, Analytics
 - Sistemas de Web Scrapping
- Los datos son muchos
- Los datos están sucios
 - Inconsistencias
 - Errores en la carga, escritura
 - Valores faltantes
 - No respuestas totales o parciales
- Cualquier error impacta en el valor que podemos extraer de un dataset

¿Qué es un valor perdido?

- Valor del que se carece una dato válido en la variable observada
- Problema generalizado en investigaciones por encuestas
- Problema cada vez más frecuente en investigaciones que usan registros administrativos o datos de redes sociales, aplicaciones, etc.
- ¿Cómo se generan esos datos perdidos?

Procesos de generación de valores perdidos

Missing Completely at Random -MCAR-

- La probabilidad de que registro tenga un valor perdido en la variable Y no está relacionada ni con los valores de Y , ni con otros valores de la matriz de datos (X)
- Los valores perdidos son una submuestra al azar de los valores totales
- ¿Cuándo no hay MCAR?
 - 1 Si algún grupo tiene mayor probabilidad de presentar datos perdidos en la variable Y y/o
 - 2 si alguno de los valores de Y tiene mayor probabilidad de presentar datos perdidos

Procesos de generación de valores perdidos

Missing at Random -MAR- y Missing Not at Random -MNAR-

- **MAR:** La probabilidad de no respuesta en Y es independiente de los valores de Y , luego de condicionar sobre otras variables
- **MNAR:** La probabilidad de no respuesta depende tanto de variables X externas, como de los valores de la variable con datos perdidos (Y)

Procesos de generación de valores perdidos

Resumen

MCAR	
X1	Y
0	NA
0	1
0	1
1	1
1	NA
2	NA
2	1
2	1
3	1
3	1
3	NA
3	1
4	1
4	NA
4	1
4	NA
4	1
4	1

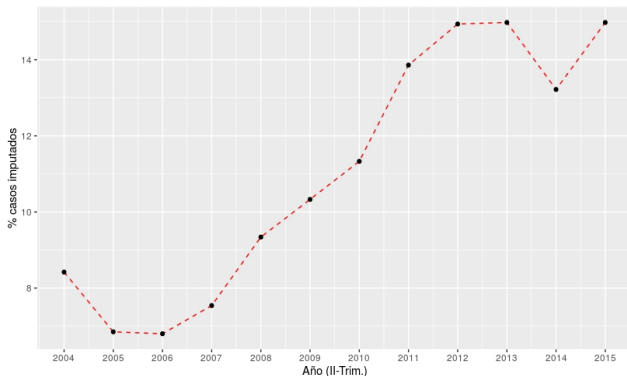
MAR	
X1	Y
0	1
0	1
0	1
1	1
1	1
2	1
2	1
2	1
3	1
3	NA
3	NA
3	1
4	1
4	NA
4	1
4	NA
4	1
4	NA

MAR	
X1	Y
0	1
0	1
0	1
1	1
1	1
2	1
2	1
2	1
3	1
3	1
3	1
3	1
4	1
4	NA
4	1
4	NA
4	NA
4	NA
4	NA

¿Por qué es importante imputar datos?

Un ejemplo: EPH

Proporción de casos imputados (sin datos en alguna variable de ingresos) en EPH. Total de aglomerados urbanos, 2003-2015 (II-Trimestre de cada año)



¿Cómo lidiar con los datos perdidos?

¿Cómo lidiar con valores perdidos?

Imputación simple

- **Excluir los casos:** se trabaja solamente con los casos completos en toda la base o solamente en las variables de estudio. Problema: se achica el dataset.
- **Reemplazar por la media o alguna otra medida:** Problema: reducción de la variabilidad de la información y se generan intervalos de confianza más estrechos de forma artificial.
- **Reponderación:** se recalculan los ponderadores de la muestra (a partir de algoritmos de reweighting) para compensar el efecto de los casos con información faltante. Problema: es incómodo trabajar con varios sets de pesos.

¿Cómo lidiar con valores perdidos?

Métodos de imputación simple - Hot Deck

- Método ampliamente usado. INDEC -hasta 2015- y Dirección de Estadística de la Ciudad para realizar imputaciones en EPH y EAH
- Reemplazar valores perdidos de un no respondente (“receptor”) con los valores observados de un respondente (“donante”) que es similar al receptor.

¿Cómo lidiar con valores perdidos?

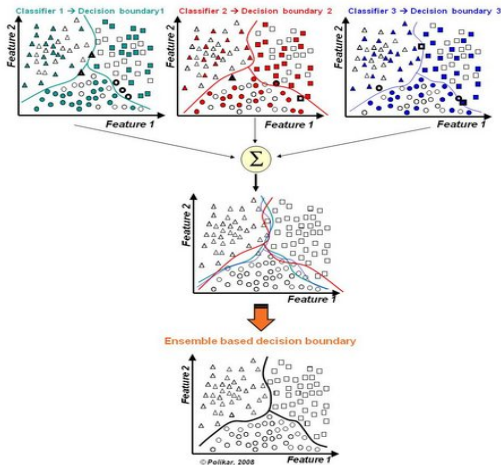
Métodos de imputación simple - Hot Deck

- **Problema 1:** selección de la " métrica" de similitud entre los casos
- **Problema 2:** selección de los donantes. El donante es seleccionado aleatoriamente de un set de potenciales donantes –hot-deck aleatorio- o bien se selecciona un solo caso donante, generalmente a partir de un algoritmo de “vecinos cercanos” usando alguna métrica -hot-deck determinístico-.

¿Cómo lidiar con valores perdidos?

Ensamble Learning

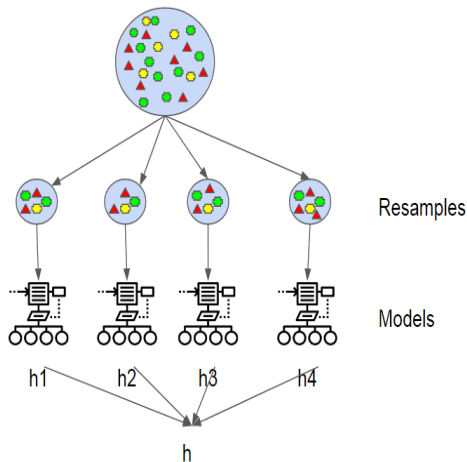
- Técnicas de aprendizaje supervisado donde se combinan varios modelos base.
- Ampliar el espacio de hipótesis posibles para mejorar la precisión predictiva del modelo combinado resultante.
- Los ensambles suelen ser mucho más precisos que los modelos base que los componen.



¿Cómo lidiar con valores perdidos?

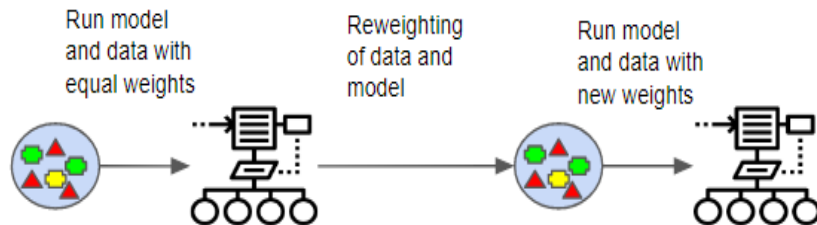
Ensamble Learning - Bagging

- Construcción de estimadores independientes -Bootstrap-
- Combinación las predicciones mediante función agregación.
- Ejemplos: Random Forest, ExtraTrees, etc.



¿Cómo lidiar con valores perdidos?

Ensamble Learning - Boosting

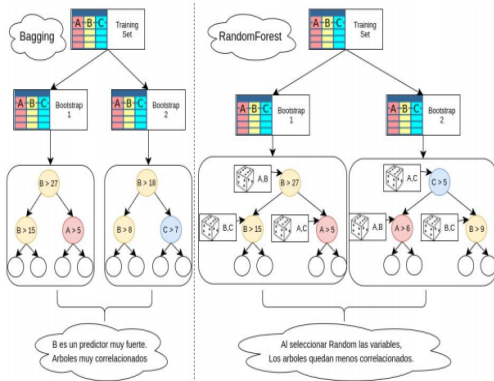


- Construcción secuencial de los estimadores
- Mayor peso en aquellos casos en los que se observa una peor performance.
- Ejemplos: AdaBoost y Gradient Tree Boosting, XGBoost.

¿Cómo lidiar con valores perdidos? - Método posible 1

Ensamble Learning - Random Forest

- Variación del algoritmo Bagging
- Modelo base: árboles de decisión
- Feature Bagging, en cada iteración y en cada split del árbol de decisión, el algoritmo selecciona aleatoriamente un subconjunto de variables predictoras



Un ejercicio de aplicación

LAB: Imputando datos perdidos con Random Forest

- Dataset: EPH 2do. trimestre de 2015
- Población: Ocupados en la semana de referencia
- Objetivo: Generar un imputador de la variable ingresos
- Variables predictoras sociodemográficas, laborales y otros ingresos
- Pipeline
 - 1 Partición Train-Test
 - 2 Sobre Train: entrenamos un clasificador Random Forest
 - 3 Sobre Train: entrenamos dos imputadores Hot-Deck
 - 4 Sobre Test: evaluamos los resultados

¿Cómo lidiar con valores perdidos? - Método posible 2

Bagging-LASSO



¿Cómo lidiar con valores perdidos? - Método posible 2

Bagging-LASSO

- Se aplica el algoritmo bagging a la imputación de ingresos laborales en la EPH del II trimestre de 2015
- En cada remuestra se estima la siguiente regresión LASSO

$$\log_{10}(y_i) = \beta_0 + \sum_{j=1}^p X_{ij}\beta_j + e_i \quad (1)$$

- Buscando minimizar la siguiente función de costo:

$$CF = RSS + \lambda \sum_{j=1}^p |\beta_j| \quad (2)$$

¿Cómo lidiar con valores perdidos? - Método posible 2

Bagging-LASSO

- ① Se crea un dataset con todos los casos completos TrS
- ② Se crea otro dataset con los casos con datos perdidos TeS
- ③ Se fija la cantidad de iteraciones rep
- ④ Para r entre 1 y rep
 - ① Se extrae una muestra bootstrap (MAS con reemplazo y $n = n^*$) de $TrSet$
 - ② En la muestra generada se estima una regresión LASSO
 - ③ Con los parámetros estimados en el paso anterior se realiza la predicción en el TeS
- ⑤ Luego de rep iteraciones se generan rep predicciones de los casos perdidos en TeS y se agregan usando la mediana

¿Cómo lidiar con valores perdidos? - Método posible 2

Bagging-LASSO

- Generación de datos perdidos de forma aleatoria
- Para evaluar performance de Bagging-LASSO y hot-deck se estima el $RMSE$ a través de k -fold cross validation, con $k = 9$
 - ① $RMSE - LASSO = 3,994\$$
 - ② $RMSE - hotdeck = 4,933\$$
- Es decir, una reducción de alrededor del 20 %
- La imputación Bagging-LASSO mejora considerablemente el error de predicción de los ingresos laborales
- Esperable que en datos perdidos “originales” (es decir, no generados artificialmente) el método consiga una mejor performance que hot-deck

- Procesos generadores de datos perdidos
- Técnicas "tradicionales"
- Ensamble Learning como alternativa
- Random Forest logra una reducción considerable en el $RMSE$ entre casos perdidos comparado a Hot Deck
- Baggin-LASSO también logra una mejora sustantiva

¿Preguntas?

german@digitalhouse.com
german.rosati@gmail.com