

## Método para predecir la probabilidad de trasplante renal para pacientes en lista de espera en Colombia

*Method for predicting the probability of kidney transplantation for patients on the waiting list in Colombia*

Cindy Zhang Gao<sup>1</sup>, Natalia Lamprea Bermudez<sup>2</sup>, Liliana López-Kleine<sup>3\*</sup>

### RESUMEN

**Introducción:** Los criterios de distribución de órganos en Colombia establecen una distribución inicial local, luego regional y por último nacional. Para diciembre de 2019 estaban en lista de espera de trasplante renal 2.822 personas en Colombia, asignadas en su mayoría a la regional Bogotá e IPS con mayores listas de espera. Esta alta concentración de pacientes podría estar generando efectos indeseados en la oportunidad que tienen los pacientes para recibir un trasplante renal. **Objetivos:** En el presente trabajo se busca estudiar, con base en datos sintéticos generados con la información disponible del INS la probabilidad de asignación de órganos identificando las variables más informativas y proponiendo un método para calcular la probabilidad de asignación para un paciente dado en lista de espera. **Material y métodos:** Se presenta el ajuste de un modelo basado en árboles de decisión que presentó una alta precisión y permite realizar la predicción de la probabilidad de obtener un órgano. **Resultados:** Se identificaron como variables más informativas el tiempo, la IPS trasplantadora y el grupo sanguíneo. Así mismo, se evidencian diferencias en

los tiempos de obtención de trasplante renal entre regionales y entre IPS trasplantadoras debido al efecto que tiene el tamaño de su lista de espera. **Conclusiones:** El método propuesto permite identificar la importancia de las variables que definen la obtención de un órgano. Finalmente, para un paciente dado, es posible estimar la probabilidad de ser clasificado en alguna de las categorías de desenlace.

**PALABRAS CLAVE:** trasplante de órganos; criterios de distribución de órganos; datos sintéticos; sistema de donación de órganos

### ABSTRACT

**Introduction:** The criteria for organ distribution in Colombia establish an initial local distribution, then regional and finally national. In December 2019, 2,822 people were on the kidney transplant waiting list in Colombia, assigned mostly to the Bogotá and IPS regions with the largest waiting lists. This high concentration of patients could be generating unwanted effects on the opportunity that patients have to receive a kidney transplant. **Objectives:** In this paper we seek to study, based on synthetic data generated with the information available from the INS,

*Correspondencia:*  
Liliana López-Kleine  
ORCID:  
0000-0001-9325-9529  
llopezk@unal.edu.co

*Financiamiento:*  
Ninguno.

*Conflicto de intereses:*  
Natalia Lamprea es paciente renal.

Recibido: 15-09-2021  
Aceptado: 10-02-2022

1) Estadística, Universidad Nacional de Colombia, Bogotá, Colombia <https://orcid.org/0000-0003-0102-1906>

2) Bióloga, Magister en Ciencias Biológicas, Especialista en Propiedad Industrial, Derecho de Autor y Nuevas Tecnologías. Universidad Nacional de Colombia, Bogotá, Colombia [natalia.lamprea@gmail.com](mailto:natalia.lamprea@gmail.com). <https://orcid.org/0000-0002-4628-1085>

3\*) Bióloga, MSc en Biometría, PhD en Estadística Aplicada, Profesora titular del Departamento de Estadística, Universidad Nacional de Colombia, Bogotá, Colombia

the probability of organ allocation, identifying the most informative variables and proposing a method to calculate the probability of allocation for a given patient on the waiting list of organs. **Material and methods:** The adjustment of a model based on decision trees is presented, which showed a high precision and allows the prediction of the probability of obtaining an organ. **Results:** Time, transplant IPS and blood group were identified as the most informative variables. Likewise, there are differences in the time it takes to obtain kidney transplants between regions and between transplanting IPS due to the effect of the size of their waiting list. **Conclusions:** The proposed method allows us to identify the importance of the variables that define obtaining an organ. Finally, for a given patient, it is possible to estimate the probability of being classified in one of the outcome categories.

**KEYWORDS:** organ transplant; organ distribution criteria; synthetic data; organ donation system

## INTRODUCCIÓN

El trasplante renal es la mejor terapia de reemplazo renal para aquellos pacientes que tienen enfermedad renal en etapa terminal.<sup>(1)</sup> La asignación de trasplante renal de donante fallecido en Colombia se realiza según el criterio de consenso de asignación establecido en el año 2018, en el cual se tienen en cuenta variables biológicas como el grupo sanguíneo, compatibilidad HLA (para los locus HLA, A, B, DR y DQ), valor calculado del panel reactivo de anticuerpos (cPRA) y diferencia de edad entre donante y receptor; variables relacionadas con el estado del receptor como es el estado compasivo y tiempo en lista de espera; y variables sociodemográficas como la ubicación geográfica.<sup>(2)</sup> Para cada una de estas variables se establece un puntaje que determina al final la asignación del órgano.

En cuanto a los criterios de distribución de órganos, se realiza como primera opción la asignación local, es decir, en las Instituciones Prestadoras de Servicios de Salud (IPS) de trasplante donde está asignado el paciente, luego la asignación regional y por último la nacional. Esta definición sigue lo establecido en el Decreto 2493/2004, art. 25. En Colombia hay seis regionales con diferente número de IPS de trasplante cada una, siendo la regional Bogo-

tá la que tiene el mayor número de IPS trasplantadoras con 13 grupos de trasplante activos.

En el 2020 la Corte Constitucional identificó una posible desigualdad en la distribución de órganos en Colombia, puesto que no se están teniendo en cuenta las diferencias en densidad poblacional y el número de pacientes que esperan trasplante en cada regional. Esto genera mayores tiempos en lista de espera, por ejemplo, para los pacientes que están ubicados en Bogotá. Por ello, la Corte le pide al Instituto Nacional de Salud (INS) que, “al momento de diseñar, modificar o desarrollar su política de distribución de órganos en general, tenga en cuenta los signos de alarma establecidos” (Sentencia T-062/20). A la fecha (septiembre 2021), el INS no ha realizado estos ajustes.

Para diciembre de 2019 estaban en lista de espera de trasplante renal 2822 personas en Colombia,<sup>(3)</sup> de las cuales 1.718 (60.8%) estaban asignadas a la regional Bogotá (Sistema de Información Instituto Nacional de Salud y Sistema de información de la Secretaría Distrital de Salud, Coordinación Regional N°. 1, Red Trasplantes 2015-2020), donde además dos de las IPS trasplantadoras de esta regional presentan el mayor tamaño en su lista de espera, con 641 pacientes en la IPS Mederi y 562 en la Fundación Cardioinfantil, lo que corresponde al 37,3% y 32,7%, respectivamente, de los pacientes de la regional Bogotá.<sup>(3)</sup> Esta alta concentración de pacientes podría estar generando efectos indeseados en la oportunidad y efectividad que tienen los pacientes para recibir un trasplante renal.

## OBJETIVOS

Este artículo tiene por objeto establecer cuáles variables entre los criterios de asignación y distribución de trasplante renal son los más determinantes para la obtención efectiva del órgano y especialmente, valorar si el tamaño de la lista de espera en la IPS de trasplante afecta la probabilidad de obtención de un riñón para el paciente asignado en ésta. Se realizó un análisis descriptivo de la información disponible del INS. Con base en la información pública que brinda el INS y otras características de la población colombiana se elaboró una base de datos sintética que se utilizó para construir árboles de decisión con el fin de predecir y cuantificar la asignación de un órgano para pacientes en lista de espera para

trasplante renal. El método propuesto muestra ser una herramienta útil que podrá ser utilizada para el caso colombiano o para otros sistemas, ya que es flexible y puede ser aplicado con la información existente. Igualmente, permite identificar la importancia de las variables que definen la asignación. Finalmente, para un paciente dado, es posible estimar la probabilidad de ser clasificado en alguna de las categorías de desenlace.

## MATERIALES Y MÉTODOS

### 1- Análisis descriptivo y obtención de las bases de datos

#### 1.1 Datos agrupados del INS

El INS publica anualmente el reporte sobre trasplantes con los datos agrupados por frecuencias y categorías. Estos reportes se pueden acceder en la página de donación de órganos y tejidos en la página del INS. Sobre esta información se realizó un análisis estadístico descriptivo univariado y multivariado utilizando los informes de los años 2016-2019. El cálculo de porcentaje de lista evacuada en un grupo de trasplante se realiza a partir del N° de pacientes en lista de espera en el año inmediatamente anterior y el N° de trasplantes realizados en el año.

#### 1.2 Obtención de los datos desagregados sintéticos

Para los análisis posteriores se utilizaron datos sintéticos que se obtuvieron a partir de la combinación de los datos reales obtenidos del informe anual del INS<sup>(4-5)</sup> y una simulación de la base de datos por individuos, ya que con los datos agrupados no era posible evaluar y aplicar el método propuesto. Se crearon dos tablas: una con individuos trasplantados y otra con individuos en lista de espera, ambas con las mismas variables. La información de estas tablas se usó para poder generar los datos sintéticos desagregados a través de distribuciones y proporciones con las características de los individuos. Se utilizaron las funciones *sample* (genera valores de unos valores conocidos) y *runif* (genera datos aleatorios de una distribución uniforme con rangos especificados) de R<sup>(6)</sup> ingresando como parámetros las frecuencias reportadas en los datos disponibles (**Tabla 1**). En un segundo tiempo se combinaron las dos tablas (trasplantados y lista de espera) y los individuos de la tabla se clasificaron en cuatro categorías: VIVO, CADAVERÍCO (para identificar tipo de trasplante de los pacientes trasplantados) y MUERE, NO MUERE (para el estado en que continúan los pacientes en lista de espera).

**Tabla 1.** Información sobre tiempo en lista de espera en los datos del informe del 2018<sup>(4)</sup>

Regional	N° días lista espera (donante cadavérico)			N° días lista espera (donante vivo)		
	Promedio	Mínimo	Máximo	Promedio	Mínimo	Máximo
1 Bogotá	935	3	3701	389	1	2583
2 Antioquia	373	2	1980	362	1	838
3 Valle	332	2	2173	263	1	1126
4 Santander	178	7	1562	1	1	1
5 Atlántico	805	6	2367	316	1	1946
6 Huila	641	83	2220	-	-	-

Los detalles de simulación para las 11 variables usadas se describen a continuación:

#### 1.2.1 Generación de variables demográficas:

Se utilizó la tabla de variables demográficas presentada en el informe del INS del 2018. Mediante la función *sample* de R, agregándole la probabilidad como parámetro a cada función de la variable se generaron los datos (**Tabla 2**).

Sexo. Para donante vivo y cadavérico se utilizan las probabilidades dadas por el informe del INS (**Tabla 2**), pero como no existe información sobre

los pacientes sin donante se simula utilizando la proporción de hombres y mujeres en Colombia según reporta en Departamento Administrativo Nacional de Estadística<sup>(5)</sup> en 2018, 49% y 51%, respectivamente.

Edad. La edad se encuentra categorizada como se muestra en la **Tabla 2** para donante vivo y cadavérico, pero para receptores sin donante se calcula una proporción a partir de la suma de donante cadavérico y vivo utilizando la función *sample* de R.

**Tabla 2.** Resumen de la información utilizada para la generación de la base de datos sintética en donde hay información o del informe de trasplantes del INS y de otras fuentes

Variable	Donante Vivo y cadavérico			Paciente en LE y Fallecido			
	Fuente de Información	Categorías	Frecuencia	Fuente de Información	Categorías	Frecuencia	
Edad	Informe Ejecutivo INS 2018 <sup>(4)</sup>	<1 año	0 %	0,30%	Informe Dane 2018 <sup>(5)</sup>	<1año	0,00 %
		1 a 5 años	0 %	0,70%		1 a 5 años	0,50 %
		6 a 11 años	2,80%	1,70%		6 a 11 años	1,90 %
		12 a 17 años	6,40%	3,30%		12 a 17 años	3,80 %
		18 a 28 años	27,00%	12,00%		18 a 28 años	15,00 %
		29 a 59 años	55,30%	63,20%		29 a 59 años	62,00 %
		> de 60 años	8,50%	18,80%		> de 60 años	17,00 %
Sexo	Informe Ejecutivo INS 2018	Femenino	46,80%	40,50%	Informe Dane 2018 <sup>(5)</sup>	Femenino	51,20%
		Masculino	53,20%	59,50%		Masculino	48,80%
ABO	Informe Ejecutivo INS 2018 <sup>(4)</sup>	O+	61,00%	56 %	Frecuencia de grupos sanguíneos ABO y Rh en donantes del banco de sangre del Hospital Pablo Tobón Uribe <sup>(6)</sup>	O+	52,10%
		O-	3,50%	3 %		O-	7,00%
		A+	24,80%	26 %		A+	28,00%
		A-	0,00%	1 %		A-	2,50%
		B+	7,80%	10 %		B+	6,70%
		B-	0,70%	0 %		B-	0,80%
		AB+	1,40%	2 %		AB+	1,70%
		AB-	0,00%	0 %		AB-	0,30%
		sin dato	<b>0,70%</b>	<b>1 %</b>		sin dato	<b>9,00%</b>
HLA*	The allele Frequency <sup>(8)</sup>	A / 24:02 -02:01-03:01		19.5% -18.72%- 7.6%			
		B / 40:02-35:43-51:01		10.17%-8%-6.7%			
		DQ / 03:02 -03:01 -02:01		21.8% -18.72% -14.91%			
		DRB / 04:07-07:01 -08:02		10.8% -9.13% -6.7%			
IPS**	Base de datos sintética	M	15,72%	Informe Ejecutivo INS 2018 <sup>(4)</sup>	M	13.87%	
		P	8,55 %		P	7.98%	
		V	13,99%		V	<b>13.18%</b>	
		F	6,13 %		F	5.32%	
		A	6,01 %		A	6.01%	
		E	3,01 %		E	3.12%	
Regional***	Informe Ejecutivo INS 2018 <sup>(4)</sup>	Regional N° 1 - Bogotá		62,34 %			
		Regional N° 2 - Antioquia		13,43 %			
		Regional N° 3 - Valle		13,43 %			
		Regional N° 4 - Santander		3,40 %			
		Regional N°5 - Atlántico		6,63 %			
		Regional N° 6 - Huila		0,78 %			
RECIBE	Informe Ejecutivo INS 2018 <sup>(4)</sup>	VIVO		18,49 %			
		CADAVERICO		81,50 %			
		MUERE		5,60 %			
		NO MUERE		94,33 %			

\* 3 primeros alelos más frecuentes por categoría de HLA. \*\* Primeras IPS más frecuentes en el sistema de trasplante de cada regional. \*\*\* Regional N° 1: Bogotá, D. C., Cundinamarca, Tolima, Boyacá, Casanare, Meta, Caquetá, Vichada, Vaupés, Guaviare, Guainía, Putumayo y Amazonas. • Regional N° 2: Antioquia, San Andrés y Providencia, Chocó, Córdoba y Caldas. • Regional N° 3: Valle, Risaralda, Quindío, Cauca y Nariño. • Regional N° 4: Santander, Norte de Santander, Cesar y Arauca. • Regional N° 5: Atlántico, Bolívar, Magdalena, Guajira y Sucre. • Regional No.6: Huila, Mederi (M), Valle de Lili (V), Paul de Rionegro(P), Fundación cardiovascular (F), C. Asunción (A), ESE.UNI Hernando Moncaleano Perdomo (E)

### 1.2.2 Variables clínicas

Estas variables fueron generadas usando igualmente la función *sample* de R asignándole una probabilidad de ocurrencia a cada categoría de la variable:

- Tipo de Rh: se utilizaron las frecuencias dadas en las tablas del informe de trasplante renal del 2018 para donante vivo y cadavérico y para los pacientes que no recibieron órgano se utiliza la probabilidad del tipo de sangre según la proporción en Colombia.<sup>(7)</sup>

- HLA: se consideraron 5 variantes en la tabla para el HLA, las cuales son combinaciones de los tipos de alelo. Al no poseer esta información ni de individuos trasplantados ni en lista de espera en los informes del INS, se usaron frecuencias reportadas por The Allele Frequency website<sup>(8)</sup> sobre una muestra de 254 individuos colombianos. Los tipos de HLA incluidos fueron, por ser los de mayor frecuencia en esta muestra, los siguientes: HLA\_A, HLA\_B, HLA\_DQ, HLA\_DRB (**Tabla 2**). Otros alelos, como HLA\_C o la variable cPRA (probabilidad reactiva al antígeno) no fueron incluidos por falta de información tanto de los informes del INS (cPRA solamente existe para los pacientes en lista de espera) y porque tampoco se dispone de una referencia para la población colombiana en general. Sin embargo, pueden incluirse en la generación de datos sintéticos o en la predicción si la información estuviera disponible.

### 1.2.3 Variables adicionales del INS

Se trata de variables que incluye el INS y están reportadas en el informe del 2018:

- IPS. Los trasplantes renales fueron realizados por 31 IPS, según el informe y así se obtienen las frecuencias de cada una para generar los datos simulados. Se utilizó la función *sample* de R, anidado a un bucle dentro de cada regional para evitar inconsistencias.

- TIEMPO EN LISTA DE ESPERA (TIEMPO). El Tiempo en lista de espera de todos los pacientes en la tabla de datos sintética fue obtenido a través de la función *runif* en R, donde se daba el máximo y el mínimo según las tablas del informe anual de trasplante renal.

- REGIONAL. La regional es generada a partir de las tablas presentadas por el INS utilizando la función *sample* de R.

- RECIBE. Esta es la variable respuesta que se construye teniendo en cuenta la cantidad de pacientes que reciben el trasplante (SI) ,

los que fallecen en lista de espera (MUERTO) y los restantes que aún continúan en lista de espera (NO).

## 2-Estimación de la probabilidad de asignación de un órgano con base en árboles de decisión

### 2.1 Árboles de decisión

El método de árbol de decisión como clasificador pertenece al grupo de algoritmos supervisados de Machine Learning. El objetivo en este caso es predecir una variable o etiqueta para los individuos a partir de un conjunto de variables predictoras. La construcción de un árbol de decisión consiste en un conjunto que se forma a partir de un nodo raíz pasando por unos nodos internos para llegar a su respuesta en un nodo hoja que es la variable que se desea predecir e indica en qué categoría de ese desenlace queda clasificado cada individuo según las ramas por las que va siendo clasificado según los nodos el árbol. El árbol se construye con base en un conjunto de aprendizaje (algunos individuos) y luego se aplica a los individuos que se desean clasificar. El método utilizado acá fue el presentado por Tangirala (2020).<sup>(9)</sup>

Sea  $L^L$  una muestra de aprendizaje

$$L = \{(x_1, c_1), (x_2, c_2), \dots, (x_p, c_p)\} \quad L = \{(x_1, c_1), (x_2, c_2), \dots, (x_p, c_p)\}$$

donde  $x_i$  representan las variables de entrada de los datos y  $c_i$  cada una de las clases de los datos en la muestra, luego las tuplas  $(x_p, c_p)$  representan las variables pertenecientes a cierta clase. Además, tenemos que  $p_i = \frac{c_i}{L} p_i = \frac{c_i}{L}$  como la probabilidad de que una tupla pertenezca a una de las clases. Con esto, definimos también la entropía como:

$$\text{Entropy}(L) = (C1/L) \log_2 (C1/L) + (C2/L) \log_2 (C2/L) + \dots + (Cj/L) \log_2 (Cj/L)$$

$$\text{Entropy}(L) = (C1/L) \log_2 (C1/L) + (C2/L) \log_2 (C2/L) + \dots + (Cj/L) \log_2 (Cj/L)$$

$$\text{Entropy}(L) = p1 \log(p1) + p2 \log(p2) + \dots + pj \log(pj) \quad \text{Entropy}(L) = p1 \log(p1) + p2 \log(p2) + \dots + pj \log(pj)$$

$$\text{Entropy}(L) = - \sum_{i=1}^j 1p \log(p1)$$

Luego, la escogencia del nodo raíz según la **Figura 1** y sus divisiones debe hacerse estratégicamente de tal forma que el árbol sea preciso en cada fase, para esto se utilizan las dos métricas más importantes para este método:

- Impureza de Gini. Los datos presentan una pureza que en términos generales es su homogeneidad, entonces para una rama cualesquiera dos individuos que escojan en un nodo interno debe tener probabilidad 1 de ser de

la misma categoría. Los valores para cada una de las categorías se calculan así:

$$GINI(L)=1-\sum_{i=1}^j p_i^2 \quad GINI(L)=1-\sum_{i=1}^j p_i^2$$

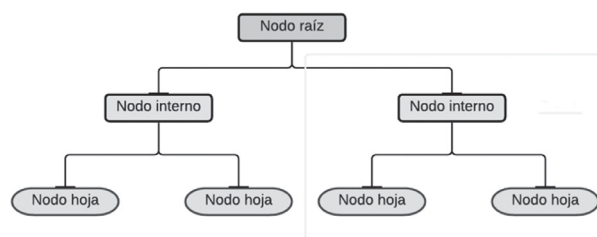
- Ganancia de información. Un nodo poco impuro tendrá menor información y la forma de medir la cantidad de información en un nodo es por medio de la entropía, donde una muestra completamente homogénea será de 0, sin embargo, para el cálculo de la Ganancia de información es necesario tener en cuenta la categoría interés y su complemento:

$$IG(L, f) = Entropy(L) - \sum_{v=1}^V \frac{|L^v|}{|L|} Entropy(L^v)$$

$$IG(L, f) = Entropy(L) - \sum_{v=1}^V \frac{|L^v|}{|L|} Entropy(L^v)$$

donde  $L^v$  representa cada muestra en  $x_i, x_i$  de cada categoría posible en la variable respuesta.

**Figura 1.** Esquema de un árbol de decisión general



Al igual que un árbol biológico, para mejorar su crecimiento requiere ser intervenido, en este caso el modelo necesitará algunos métodos que permitirán mejorarlo como la escogencia de un máximo de hojas, o un mínimo de muestras, la poda, algunos métodos de ensamble como los son los Bosques aleatorios, entre otros.<sup>(9)</sup> Estos fueron usados por la librería *Scikit-Learn* de Python y algunas librerías para procesamiento de los datos como *numpy*, *pandas*.<sup>(13)</sup>

**2.2 Validación cruzada (k-fold)**

Sobre el modelo ajustado se realiza una validación cruzada sobre una muestra de individuos para estimar los valores de exactitud para cada modelo. La técnica de validación cruzada es bastante práctica al momento de evaluar algún método de aprendizaje automático. Este consiste en dividir la muestra en k partes o pliegues (fold) para evaluar k veces el modelo siguiendo la misma metodología sobre una cantidad de los datos en entrenamiento y otro en validación. Esta etapa se hace con el fin de

evitar que los modelos se sobreajusten y dependan de una muestra específica sin que puedan ser útiles en datos diferentes <sup>(10)</sup>. Se utilizaron 5 folds y la función *cross\_val\_score* de la librería *Scikit-Learn* de Python<sup>(12)</sup>, además, se optimizaron los hiperparámetros: profundidad mínima,<sup>(6)</sup> muestras máximas y mínimas por cada hoja.<sup>(10)</sup>

**2.2 Bootstrap**

La técnica de Bootstrap es un método de remuestreo con reemplazo basado en una simulación de datos con el fin de conocer el valor real de un estadístico, para nuestro caso la proporción de pacientes pertenecientes a cada clase, donde se utilizan submuestras con reemplazo de la misma población generando intervalos de confianza y la estimación puntual de un valor.<sup>(11)</sup> Se utilizó la librería *boot* de R<sup>(14)</sup> para la estimación de las proporciones de cada clase y las funciones de base de R para manejo procesar los resultados.

**2.3 Métricas para evaluar la calidad del modelo**

Dentro de los modelos de clasificación se utilizan una técnica de etiquetado, donde se trabaja con un conjunto de entrenamiento del modelo, y validación el cual es utilizado para predecir las etiquetas de la clase de datos nuevos. Se trabaja con una matriz de confusión de la siguiente manera: **(Tabla 3)**

**Tabla 3.** Matriz de confusión sobre la predicción con la base de datos sintética con respecto a la base actual

		Actual	
		Positivo	Negativo
Predicción de la clase	Positivo	Verdadero Positivo (VP)	Falso Positivo (FP)
	Negativo	Falso Negativo (FN)	Verdadero Negativo (VN)

Estas matrices de confusión son presentadas en el material suplementario tanto para el árbol con todos los datos agrupados, como para cada árbol de las regionales.

Las métricas más utilizadas para este tipo de metodología son:

- *Accuracy*. Es la relación entre el número de predicciones correctas y el número total de predicciones realizadas. También es posible entenderla como la velocidad con que el modelo

hace predicciones.

$$accuracy = \frac{VP + VN}{VP + FP + FN + VN}$$

- *Precisión*. Proporción de valores de la clase mayoritaria correctamente clasificada (Verdaderos Positivos) dividida por la suma de los valores de la clase mayoritaria correctamente clasificados (verdaderos positivos) y los valores de la mayoría clasificados incorrectamente (falso positivo).

$$precision = \frac{VP}{VP+FP}$$

- *Recall* (exhaustividad). Proporción de valores correctamente clasificados que sean de la clase mayoritaria (verdaderos positivos) sobre los que realmente estaban en esa clase como la suma de los valores en la clase mayoritaria verdaderamente clasificados (verdaderos positivos) y los valores de la clase minoritaria clasificados incorrectamente (verdaderos negativos).

$$recall = \frac{VP}{VP+FN}$$

- *F1-Score*. Es la combinación de la precisión y el *recall* calculado como una media armónica, esta métrica es interesante en los clasificadores porque tiene ambas métricas de precisión y *recall* con la

misma importancia.

$$F1 = 2 * \frac{precision * recall}{precision + recall}$$

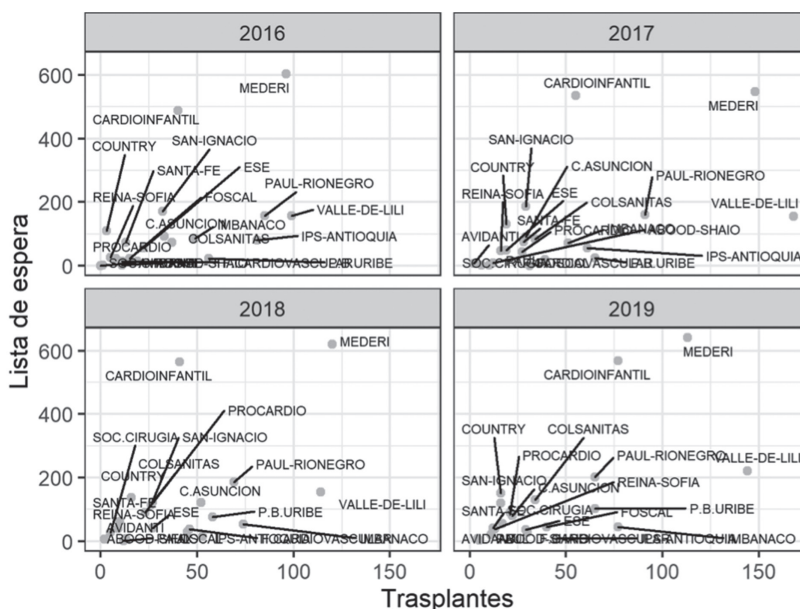
Se ha puesto a disposición todo el código utilizado en el siguiente repositorio: [https://github.com/czhangg17/articulo\\_trasplantes](https://github.com/czhangg17/articulo_trasplantes).

## RESULTADOS

### 1) Descripción de los datos a partir de informes del INS

En la **Figura 2** muestra que durante los años 2016 a 2019 dos IPS, Cardioinfantil y Mederi, mantienen un número de pacientes en lista de espera entre 500 y 600, sin lograr disminuir el tamaño de su lista a lo largo de los años; respecto al número de trasplantes por año, aunque éstos han aumentado con los años, no logran alcanzar los 100 para Cardio infantil o 150 para Mederi. En un rango medio se encuentran IPS como Hospital San Ignacio, San Vicente de Paul- Rionegro, y Valle de Lili que mantienen su lista de espera alrededor de los 200 pacientes, con un número de trasplantes entre 50 y 140. Las demás IPS tienden a mantener el número de pacientes en lista por debajo de los 100 pacientes, con un número de trasplantes por debajo de 50 al año.

**Figura 2.** Diagramas de dispersión del número de trasplantes vs el número de pacientes en lista de espera para los años 2016 a 2019 con base en datos publicados por los informes de trasplantes del INS para 22 IPS



Así, la IPS Mederi tiene el mayor número de pacientes en lista de espera y en promedio con 117 trasplantes por año, logra evacuar el 21,7%

de ésta lista; mientras que la Fundación Valle de Lili encabeza el número de trasplantes renales anualmente, en promedio con 130,5 trasplantes por

año, logra evacuar el 84,5 % de su lista de espera.

Lo observado en estos diagramas de dispersión indica que hay diferencias importantes entre IPS. Especialmente, en lo referente a la proporción número de trasplantes realizados al año y número de pacientes en lista de espera, puesto que las IPS en el grupo medio y en el menor tiene proporciones entre 1:2 a 1:4, es decir, hay un paciente trasplantado por cada dos a cuatro pacientes en lista de espera; mientras que las dos IPS del grupo superior tienen una proporción de 1:5 a 1:10, lo que significa que por cada paciente trasplantado hay entre 5 y 10 pacientes en lista de espera.

### 2. Descripción de los datos desagregados sintéticos

Una vez generada la tabla de datos desagregados

sintéticos se construye una estimación de las categorías de la variable de interés realizando un *bootstrap* (remuestreo) sobre las mismas para obtener la proporción estimada en cada categoría (**Tabla 4**). Esta proporción indica la probabilidad de un paciente, con base en la información generada, de estar en alguna de las cuatro categorías de la variable de interés. Lo anterior nos indica que a priori el 70,5% de las personas no mueren en la lista de espera, mientras que se estima que un total un 20,5% recibe un trasplante de un donante cadavérico y el 0,046% lo recibe de un donante vivo. Lo cual significa que aproximadamente el 79% de los pacientes en lista de espera no reciben órgano. Estas estimaciones son para el momento presente en el que se estima la probabilidad. Ya que los datos disponibles son para el año 2018, esta estimación aplicaría para ese año.

**Tabla 4.** Estimación de la proporción de pacientes para cada una de las cuatro categorías de la variable desenlace utilizando el modelo

Clase	Proporción	I.C 95%	
		Inferior	Superior
Cadavérico	0.2050	0.1913	0.2182
Muere	0.0434	0.0366	0.0500
No muere	0.705	0.6908	0.7196
Vivo	0.046	0.0398	0.0538

### 3. Árboles de decisión

Al construir un árbol de decisión para la clasificación de los pacientes en las cuatro categorías de desenlace, se encontró un buen ajuste

para pacientes pertenecientes a las regionales 1, 2 y 3 debido a la cantidad de información disponible para estas regionales (**Tabla 5**).

**Tabla 5.** Métricas de calidad de precisión, exactitud, recall y f1-score sobre los modelos de árboles de decisión construidos

Modelo	CADAVÉRICO			MUERE			NO MUERE			VIVO			Exactitud (Accuracy)
	Precisión	Recall	F1-Score	Precisión	Recall	F1-Score	Precisión	Recall	F1-Score	Precisión	Recall	F1-Score	
Mod.opt	0.82	0.8	0.81	0	0	0	0.9	0.99	0.94	0	0	0	0.88
Mod.Reg1	0.82	0.79	0.81	0	0	0	0.92	0.99	0.96	0	0	0	0.91
Mod.Reg2	0.77	0.83	0.8	0	0	0	0.76	0.89	0.82	0	0	0	0.76
Mod.Reg3	0.65	0.58	0.61	0	0	0	0.81	0.98	0.88	0.12	0.08	0.1	0.71
Mod.Reg4	0.89	1	0.94	-	-	-	1	0.67	0.8	1	1	1	0.92
Mod.Reg5	0.5	0.5	0.5	0	0	0	0.85	0.94	0.89	0	0	0	0.8
Mod.Reg6	0.67	1	0.8	0	0	0	-	-	-	0	0	0	0.67

Mod.opt: Modelo completo para todas las regionales. Mod.Reg son los modelos para cada regional: 1 Bogotá, 2 Antioquia, 3 Valle del Cauca, 4 Santander, 5 Atlántico, 6 Huila.



Se construye un árbol tomando todos los datos y optimizando la cantidad máxima de hojas y profundidad del árbol. Se obtiene una exactitud del 88% en su clasificación tomando un 80% de

la tabla para entrenar el modelo y otro 20% para validarlo (Mod. opt) con un resultado sobre la clasificación de las categorías en una matriz de confusión (**Tabla suplementaria 2**).

**Tabla suplementaria 2.** Matriz de confusión para el Mod.Opt-Modelo completo para todas las regionales

Mod.opt	Valores reales			
	CADAVÉRICO	MUERE	NO MUERE	VIVO
Valores predichos				
CADAVÉRICO	153	0	34	5
MUERE	0	0	36	0
NO MUERE	5	0	756	0
VIVO	29	0	15	0

Adicionalmente, se realiza el mismo modelo para cada una de las regionales para lograr tener un modelo independiente de cada regional y poder estimar sin los pesos de otros pacientes pertenecientes a regionales diferentes. Las métricas de los modelos son buenas. Sin embargo, para los modelos 5 y 6 no es posible tener exactitud aceptable sobre la clasificación debido a la calidad de los datos sintéticos que fueron generados con muy poca información (muy bajo número de

pacientes), de hecho, para el modelo 6 solamente presenta dos clases a clasificar (CADAVÉRICO y VIVO) como se observa en la matriz de confusión (**Tabla suplementaria 8**). De la misma forma, tenemos el modelo sobre la regional 5, que para la época tenía solo una entidad trasplantadora. Sin embargo, los resultados del clasificador son presentados en su matriz de confusión (**Tabla suplementaria 7**) sin pacientes que pertenecen a la clase que mueren.

**Tabla suplementaria 7.** Matriz de confusión para el Mod.Reg5-Modelo de la regional Atlántico

Mod.Reg5	Valores reales			
	CADAVÉRICO	MUERE	NO MUERE	VIVO
Valores predichos				
CADAVÉRICO	3	0	3	0
MUERE	0	0	2	0
NO MUERE	2	0	34	0
VIVO	1	0	1	0

**Tabla suplementaria 8.** Matriz de confusión para el Mod.Reg6-Modelo de la regional Huila

Mod.Reg6	Valores reales	
	CADAVÉRICO	VIVO
Valores predichos		
CADAVÉRICO	4	0
VIVO	2	0

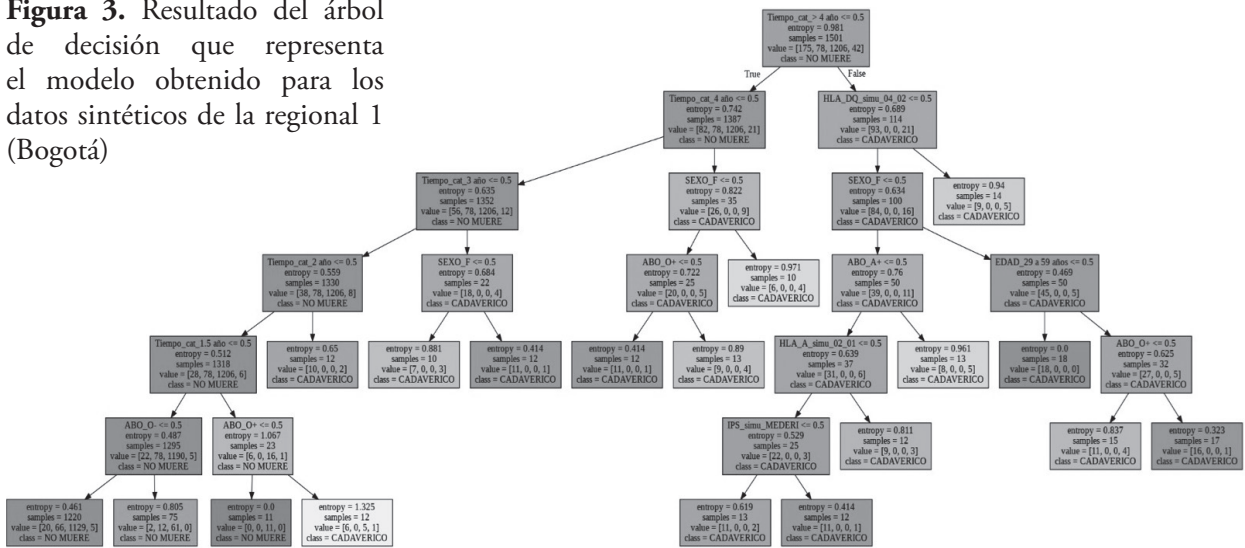
La **Figura 3** muestra el árbol de decisión obtenido para la regional 1 en el cual se puede apreciar que la variable que mayor cantidad de información proporciona para predecir el desenlace es el Tiempo (tiempo en días en la lista de espera),

la **Tabla suplementaria 3** muestra que la mayor frecuencia se encuentra en la diagonal sobre la clase de pacientes que quedan esperando en lista de espera (NO MUERE). El umbral de “mayor a 4 años” es el nodo principal del árbol y de varios

subárboles que se generan con nodos definidos por puntos de corte en otras variables como HLA, sexo o tipo de sangre. Por otro lado, algunos de los valores del HLA bifurcan entre la clasificación de recibir un órgano de tipo cadavérico. Igualmente, el árbol ilustra el camino de los pacientes que son más propensos a seguir en la lista de espera (**Figura 3**). La única IPS que figura en el modelo como importante para la predicción es la IPS

Mederi, la cual bifurca al final del árbol con un umbral de corte de 0.5. Es posible afirmar que esta bifurcación clasifica como donante de tipo cadavérico en cualquier caso para esta IPS. Por lo tanto, el modelo ajustado para este caso indica que es más probable que una persona A+, femenino, HLA\_A\_02\_01 tenga un desenlace de recibir un órgano de un donante cadavérico, si además tiene menos de 4 años en lista (que es su nodo principal).

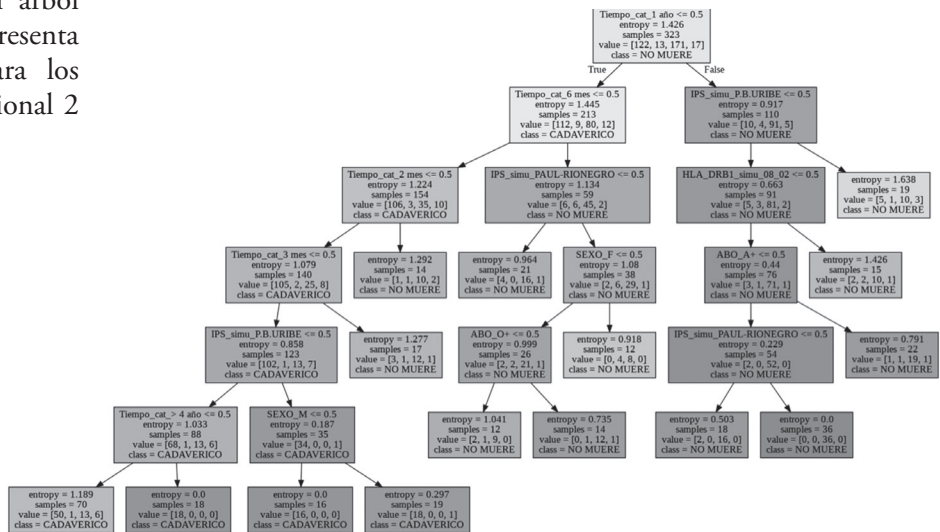
**Figura 3.** Resultado del árbol de decisión que representa el modelo obtenido para los datos sintéticos de la regional 1 (Bogotá)



El valor de entropy indica la ganancia de información aportada por ese nodo del árbol (ver sección metodología). Samples indica el número de individuos que se usaron para determinar el valor de la variable de cada nodo, el número de

individuos perteneciente a cada una de las cuatro categorías está indicado entre paréntesis cuadradas en value: [Cadavérico, Muere, No Muere, Vivo]. La suma de samples de los nodos hijos es el número de muestras del nodo padre.

**Figura 4.** Resultado del árbol de decisión que representa el modelo obtenido para los datos sintéticos de la regional 2 (Antioquia)



**Tabla suplementaria 3.** Matriz de confusión para el Mod.Reg1-Modelo de la regional Bogotá

Mod.Reg1	Valores reales			
Valores predichos	CADAVERÍCO	MUERE	NO MUERE	VIVO
CADAVERÍCO	61	0	16	0
MUERE	0	0	28	0
NO MUERE	3	0	524	0
VIVO	10	0	2	0

De la misma manera, el árbol obtenido para la regional 2-Antioquia nos permite visualizar la importancia de otras variables como por ejemplo la asignación en la IPS Pablo Tobon Uribe. Esta se extiende en el subárbol de los que quedan en lista de espera partiendo de su nodo principal el tiempo en lista de espera menor a 1 año. Este modelo, a pesar de tener valores en las métricas menores al

de la regional 1, visualmente aporta información valiosa para la comprensión de las variables más informativas. Además, la **Tabla suplementaria 4** muestra la matriz de confusión del modelo 2 donde 48 de 57 pacientes en la validación del modelo son predichos correctamente en la clase para recibir el trasplante.

**Tabla suplementaria 4.** Matriz de confusión para el Mod.Reg2-Modelo de la regional Antioquia

Mod. Reg2	Valores reales			
Valores predichos	CADAVERÍCO	MUERE	NO MUERE	VIVO
CADAVERÍCO	44	0	9	0
MUERE	1	0	5	0
NO MUERE	8	0	62	0
VIVO	4	0	6	0

Los modelos para las regionales 3 (Valle del Cauca) y 4 (Santander) se encuentran en anexo (**Figuras suplementarias 1 y 2**) junto con sus matrices de confusión (**Tablas suplementarias 5 y 6**). Se confirma que la variable más importante es el tiempo. En el modelo de la regional 3 se puede apreciar que, al superar un tiempo de 4 años en lista de espera, todos los pacientes reciben

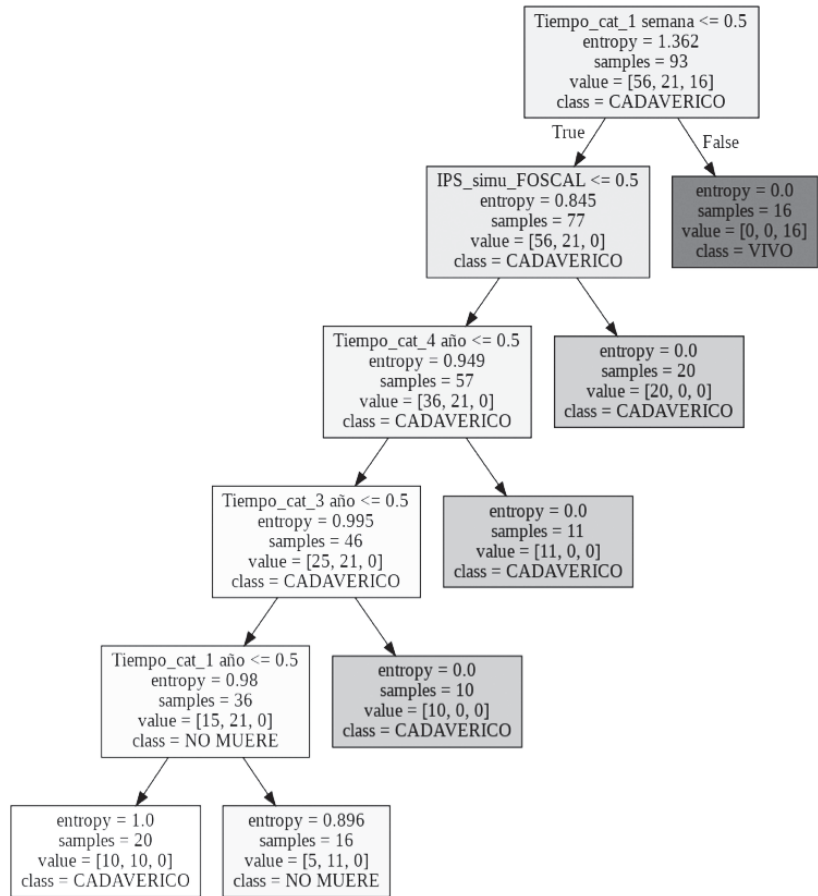
un órgano si hacen parte de la IPS Valle de Lili. También es posible observar que el modelo de la regional 4 no es informativo, ya que todos los pacientes son clasificados en la misma categoría (recibir un donante cadavérico). Es de esperarse que el modelo no sea informativo, dado que fue construido solamente con la información de 93 individuos.

**Figura suplementaria 2.** Resultado del árbol de decisión que representa el modelo obtenido para los datos sintéticos de la regional 4 (Santander)

```

Predicción de pacientes con árbol Regional 1
Los dos primeros pacientes son los mismos al los primeros en las predicciones anteriores
-----
      IPS SEXO      EDAD ABO  HLA_A  HLA_B  HLA_DQ  HLA_DRB1  \
0  CARDIOINFANTIL  M  29 a 59 años  O+  24_02  44_03  02_01  13_03
1  MEDERI          F  mayores de 60  A-  03_01  35_01  03_02  04_07
2  CARDIOINFANTIL  F  mayores de 60  A-  03_01  35_01  03_02  04_07
3  MEDERI          F  29 a 59 años  O+  03_01  35_01  03_02  04_07
-----
      TIEMPO  CADAVERICO  MUERE  NO MUERE  VIVO  CLASE
0  > 4 año  0.846154  0.000000  0.000000  0.153846  CADAVERICO
1  1 año  0.016393  0.054098  0.92541  0.004098  NO MUERE
2  1 año  0.016393  0.054098  0.92541  0.004098  NO MUERE
3  > 4 año  0.941176  0.000000  0.000000  0.058824  CADAVERICO
-----
    
```

**Figura suplementaria 1.** Resultado del árbol de decisión que representa el modelo obtenido para los datos sintéticos de la regional 3 (Valle del Cauca)



**Tabla suplementaria 5.** Matriz de confusión para el Mod.Reg3-Modelo de la regional Valle del Cauca

Mod.Reg3	Valores reales			
	CADAVERÍCO	MUERE	NO MUERE	VIVO
Valores predichos CADAVERÍCO	15	0	4	7
MUERE	0	0	3	0
NO MUERE	1	0	50	0
VIVO	7	0	5	1

**Tabla suplementaria 6.** Matriz de confusión para el Mod.Reg4-Modelo de la regional Santander

Mod.Reg4	CADAVERÍCO	NO MUERE	VIVO
CADAVERÍCO	16	0	0
NO MUERE	2	4	0
VIVO	0	0	2

**4. Predicción de la variable desenlace**

Además de obtener un árbol de decisión que permite visualizar la importancia de las variables y sus puntos de corte, es posible usar los árboles generados para predecir, según los valores de las

variables para un paciente hipotético dado, cuál es su probabilidad de ser clasificado en alguna de las cuatro categorías de la variable desenlace (Tabla 6). En caso de que los datos sean reales, es posible realizar esta predicción para cada uno de los

pacientes de la base de datos. La **Tabla 6** muestra para ciertos valores de las variables sintéticas de la base de datos, la probabilidad estimada de ser

clasificado en cada una de las categorías por los árboles construidos presentados en la sección anterior.

**Tabla 6.** probabilidad de caer en cada categoría de la variable desenlace para 9 pacientes hipotéticos predicción del modelo completo

IPS	Sexo	Edad	ABO	HLA_A	HLA_B	HLA_DQ	HLA_DRB1	Reg.	T	P Cadav.	P Muere	P No muere	P vivo
C	M	29-59	O+	24_02	44_03	02_01	13_03	1	>4	<b>0.600</b>	0.000	0.000	0.400
M	F	>60	A-	03_01	35_01	03_02	04_07	1	1	0.015	0.062	<b>0.919</b>	0.005
M	F	29-59	O+	03_01	35_01	03_02	04_07	1	>4	<b>0.941</b>	0.000	0.000	0.058
V	M	29-59	A+	02_05	35_12	03_02	08_02	3	>4	<b>1.000</b>	0.000	0.000	0.000
V	M	>60	O+	68_01	35_12	02_01	03_02	3	1	0.092	0.056	<b>0.794</b>	0.058
I	F	18-28	A+	68_02	49_01	06_02	08_02	2	1	0.092	0.056	<b>0.794</b>	0.058
I	M	>60	B-	02_01	48_01	06_04	04_01	2	>4	<b>1.000</b>	0.000	0.000	0.000
F	M	>60	O+	03_01	40_02	03_04	13_04	4	3	<b>1.000</b>	0.000	0.000	0.000
A	M	29-59	A+	02_13	35_43	02_01	01_02	5	1	0.092	0.056	<b>0.794</b>	0.058
E	M	>60	O+	03_01	35_43	03_01	14_01	6	2	<b>0.741</b>	0.000	0.000	0.260

En negrita la clasificación del paciente según la mayor probabilidad. IPS: Cardioinfantil (C), Mederi (M), Valle de Lili (V), IPS-Antioquia (I), Fundación cardiovascular (F), C. Asunción (A), ESE.UNI Hernando Moncaleano Perdomo (E). Regionales: 1 (Bogotá), 2 (Antioquia), 3 (valle del Cauca), 4 (Santander), 5 (Atlántico), 6 (Huila). T: Tiempo en años. P: probabilidad.

Aquí vemos las diferencias en los desenlaces y en los tiempos de obtener por ejemplo el trasplante cadavérico, entre regionales, así como entre IPS trasplantadoras, incluso para pacientes hipotéticos con perfiles parecidos.

En el anexo se presentan otras predicciones para los modelos de cada regional por separado (**Tabla suplementaria 1**), ya que es posible usar el modelo completo o el modelo por cada regional. Se resalta que los pacientes hipotéticos de la regional 1 (con los mismos valores en sus variables) tienen la misma clasificación en el modelo completo que en el modelo de su regional, aunque las probabilidades cambian (ver **Tabla suplementaria 1 - Pág. 238**).

**DISCUSIÓN**

Los datos descriptivos obtenidos del INS y las simulaciones realizadas dejan ver que existen diferencias en los desenlaces del paciente, especialmente en lo referente al tiempo de obtención del trasplante cadavérico entre regionales y entre IPS trasplantadoras. Coincidiendo en ambos la identificación de la Regional 1 como la regional donde más tiempo se demora obtener un

trasplante.

Las variables más determinantes para obtener el trasplante, específicamente de donante cadavérico, fueron diferentes entre las regionales, lo que deja ver cómo los factores que determinan el desenlace si están variando por condiciones como disponibilidad de donantes por regional o tamaño de la lista de espera de la IPS trasplantadora en la que está asignado el paciente. Estos resultados, corroboran las inequidades en la distribución de órganos que hay en Colombia y que fueron identificadas por la Corte Constitucional y particularmente, la necesidad de replantear la definición de distribución local, pues allí se está dando buena parte de las diferencias en la probabilidad de obtener un riñón de donante cadavérico.

Para la regional 1- Bogotá, la variable más informativa y determinante es el tiempo, siendo en su mayoría más de cuatro años en lista de espera representando el 69,9% de los pacientes de los cuales se etiquetan un 11,65% como donante de tipo cadavérico, 5,19% fallecen, 80.34% no mueren en lista de espera y 2,7% reciben órgano de

**Tabla suplementaria 1.** Predicción de la probabilidad de ser clasificado en cada una de las categorías para regionales 1, 2 y 3**Predicción de pacientes con árbol Regional 1**

Los dos primeros pacientes son los mismos al los primeros en las predicciones anteriores

	IPS	SEXO	EDAD	ABO	HLA_A	HLA_B	HLA_DQ	HLA_DRB1	\
0	CARDIOINFANTIL	M	29 a 59 años	O+	24_02	44_03	02_01	13_03	
1	MEDERI	F	mayores de 60	A-	03_01	35_01	03_02	04_07	
2	CARDIOINFANTIL	F	mayores de 60	A-	03_01	35_01	03_02	04_07	
3	MEDERI	F	29 a 59 años	O+	03_01	35_01	03_02	04_07	

	TIEMPO	CADAVERICO	MUERE	NO MUERE	VIVO	CLASE
0	> 4 año	0.846154	0.000000	0.000000	0.153846	CADAVERICO
1	1 año	0.016393	0.054098	0.92541	0.004098	NO MUERE
2	1 año	0.016393	0.054098	0.92541	0.004098	NO MUERE
3	> 4 año	0.941176	0.000000	0.000000	0.058824	CADAVERICO

**Predicción con pacientes árbol regional 2**

	IPS	SEXO	EDAD	ABO	HLA_A	HLA_B	HLA_DQ	HLA_DRB1	\
0	IPS-ANTIOQUIA	F	18 a 28 años	A+	68_02	49_01	06_02	08_02	
1	IPS-ANTIOQUIA	M	mayores de 60	B-	02_01	48_01	06_04	04_01	

	TIEMPO	CADAVERICO	MUERE	NO MUERE	VIVO	CLASE
0	1 año	0.133333	0.133333	0.666667	0.066667	NO MUERE
1	> 4 año	1.000000	0.000000	0.000000	0.000000	CADAVERICO

**Predicción con pacientes árbol regional 3**

	IPS	SEXO	EDAD	ABO	HLA_A	HLA_B	HLA_DQ	HLA_DRB1	\
0	VALLE-DE-LILI	M	29 a 59 años	A+	02_05	35_12	03_02	08_02	
1	VALLE-DE-LILI	M	mayores de 60	O+	68_01	35_12	02_01	03_02	

	TIEMPO	CADAVERICO	MUERE	NO MUERE	VIVO	CLASE
0	> 4 año	1.000000	0.000000	0.000000	0.000000	CADAVERICO
1	1 año	0.061674	0.035242	0.881057	0.022026	NO MUERE

donante vivo. Además, se observa que el desenlace de obtener trasplante de donante cadavérico para el grupo sanguíneo más común, que es el O+, aquí puede tomar más de cuatro años, lo que podría estar explicado por el efecto de competencia o alta-selectividad que surge al tener tan concentrados los pacientes en dos grupos de trasplante. Para la regional 2- Antioquia, el tiempo de menos de un año es la variable más informativa, seguido de menos de seis meses en lista de espera y logra el desenlace de trasplante de donante cadavérico. La asignación a la IPS Pablo Tobón Uribe, también se resalta como informativa para tener este desenlace. Para la regional 3-Valle, la variable más

determinante es el tiempo.

El método de predicción propuesto permite tener una exactitud superior al 90% en cuanto a la relevancia de las variables y las predicciones que se pueden hacer con este. Esta exactitud puede variar según la calidad y cantidad de los datos. La utilización de los datos reales de los pacientes en lista de espera permitiría conocer y comparar los desenlaces que se tienen entre las diferentes regionales y estimar las variables que los determinan. Así mismo, ante la necesidad de realizar ajustes en los criterios de distribución de órganos, el método aquí propuesto permitiría ensayar las diferentes aproximaciones de propuestas de ajustes y determinar cuáles serían las

más adecuadas para el modelo colombiano.

Las simulaciones aquí realizadas tienen la limitación que no se pudieron tener en cuenta variables importantes para los criterios de asignación como son los valores de cPRA, o el estado clínico del paciente, o incluir variables de índole administrativa como es la EPS de afiliación del paciente, que pueden explicar las diferencias obtenidas entre los datos descriptivos y los desenlaces de las simulaciones. Para las categorías de frecuencia baja (tales como trasplante de donante vivo o muere) es difícil lograr un ajuste adecuado. Así mismo, para el desenlace “muere” es necesario tener la información que permita estimar los factores que determinan realmente ese desenlace, para el presente estudio no fue posible tener esta información ni hacer esta estimación de manera real, ya que no se disponía de los datos necesarios.

Por último, queremos unirnos al llamado que han realizado otros investigadores,<sup>(12)</sup> para que los datos sobre la actividad de donación y trasplante sean abiertos, de manera que cualquier grupo pueda trabajar con ellos, esto le aporta mayor transparencia, veeduría y seguimiento al sistema. Además, permitiría tener una mayor participación de la academia, las autoridades administrativas y de los mismos pacientes, en la formulación de políticas que busquen optimizar la red de donación y trasplante en Colombia.

## BIBLIOGRAFÍA

- 1) Tonelli M, Wiebe N, Knoll G, et al. Systematic review: kidney transplantation compared with dialysis in clinically relevant outcomes. *Am J Transplant.* 2011; 11(10):2093-2109.
- 2) Instituto Nacional de Salud. 2018. Criterios de Asignación para Trasplante Renal en Colombia. Fecha de consulta: Marzo 2021. Disponible: <http://www.ins.gov.co/Direcciones/RedesSaludPublica/DonacionOrganosYTejidos/Paginas/default.aspx>
- 3) Informe Nacional de Salud. Red Nacional de Donación y Trasplantes 2019. Informe Ejecutivo. Disponible: <https://www.ins.gov.co/BibliotecaDigital/informe-ejecutivo-red-donacion-y-trasplantes-2019.pdf>
- 4) Instituto Nacional de Salud. Informe anual red de donación y trasplantes. Colombia, 2018. Disponible: <https://www.ins.gov.co/BibliotecaDigital/informe-anual-red-de-donacion-trasplantes-2018.pdf>
- 5) Departamento Nacional de Estadística. Población de Colombia es de 48,2 millones de habitantes, según el DANE. Disponible en: <https://id.presidencia.gov.co/Paginas/prensa/2019/190704-Poblacion-de-Colombia-es-de-48-2-millones-habitantes-segun-DANE.aspx>
- 6) R Core Team (2020). R: A language and environment for statistical computing. R Foundation for Statistical Computing, Vienna, Austria. Disponible: <https://www.R-project.org/>.
- 7) Beltrán M, Ayala M, Jara J. Frecuencia de grupos sanguíneos y factor Rh en donantes de sangre, Colombia, 1996. *Biomédica.* 1999; 19(1): 39-44.
- 8) Allele Frequency Net Database . Disponible: <http://www.allelefrequencies.net/>
- 9) Tangirala S. Evaluating the Impact of GINI Index and Information Gain on Classification using Decision Tree Classifier Algorithm\*. *Int J Adv Comput Sci Appl.* 2020;11(2): 612-19.
- 10) Elkan Ch. Evaluating Classifiers. University of California, San Diego, 18 de enero de 2011. Disponible en: <https://web.archive.org/web/20111218192652/http://cseweb.ucsd.edu/~elkan/250B/classifiereval.pdf>
- 11) Efron B. Bootstrap Methods: Another Look at the Jackknife. *The Annals of Statistics.* 1979; 7(1):1-26.
- 12) Lima B. A call for open data of renal transplantation in Portugal. *Port J Nephrol Hypert.* 2017; 31(3): 155-7.
- 13) F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, J. Vanderplas. Scikit-learn: machine learning in Python *J. Mach. Learn. Res.*, 12 (2011), pp. 2825-2830.
- 14) Canty A, Ripley BD (2021). *boot: Bootstrap R (S-Plus) Functions.* R package version 1.3-28.