

Aplicación de Ciencia de Datos en la Optimización de modelos predictivos para la estimación de mortalidad en Ictus Isquémico

Luis Téllez Ramírez
Jesús Martínez Gómez
Juan Manuel García Torrecillas

Universidad de Castilla La Mancha

Defensa del TFM 20/21

Índice

1 Introducción

- ¿Por qué estudiar este problema?
- CMBD: Conjunto Mínimo Básico de Datos
- Objetivos del proyecto

2 Análisis Exploratorio de Datos

- 2 escenarios
- Recorrido por las variables más importantes
- Prevalencia de los predictores
- Discusión de la métrica y división muestra

3 Metodología

- Datos desbalanceados
- Algoritmo Genético de Selección de Predictores

4 Escenarios

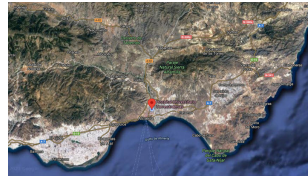
- Escenario A
 - Modelos Utilizados
- Escenario B
 - Modelos Utilizados
- Conclusiones y Trabajo Futuro

Antes de empezar

- TFM completamente distinto.



(a)



(b)

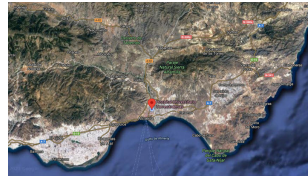
- Problema abierto a día de hoy (abierto a sugerencias).
- Grupo de investigación (Juan Manuel García Torrecillas).
- Usar las técnicas obtenidas a lo largo del máster.
- Aportar otro enfoque nuevo y diferente.

Antes de empezar

- TFM completamente distinto.



(c)



(d)

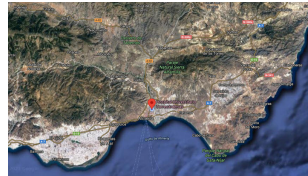
- Problema abierto a día de hoy (abierto a sugerencias).
- Grupo de investigación (Juan Manuel García Torrecillas).
- Usar las técnicas obtenidas a lo largo del máster.
- Aportar otro enfoque nuevo y diferente.

Antes de empezar

- TFM completamente distinto.



(e)



(f)

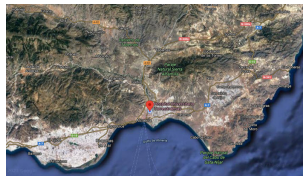
- Problema abierto a día de hoy (abierto a sugerencias).
- Grupo de investigación (Juan Manuel García Torrecillas).
- Usar las técnicas obtenidas a lo largo del máster.
- Aportar otro enfoque nuevo y diferente.

Antes de empezar

- TFM completamente distinto.



(g)



(h)

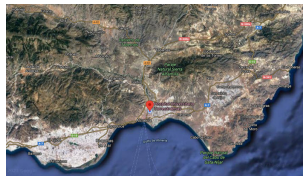
- Problema abierto a día de hoy (abierto a sugerencias).
- Grupo de investigación (Juan Manuel García Torrecillas).
- Usar las técnicas obtenidas a lo largo del máster.
- Aportar otro enfoque nuevo y diferente.

Antes de empezar

- TFM completamente distinto.



(i)



(j)

- Problema abierto a día de hoy (abierto a sugerencias).
- Grupo de investigación (Juan Manuel García Torrecillas).
- Usar las técnicas obtenidas a lo largo del máster.
- Aportar otro enfoque nuevo y diferente.

¿Por qué es un problema?

- El ictus isquémico supone la segunda causa de mortalidad en nuestro país en la población general [1].
- La primera causa de mortalidad en la mujer [1].
- A nivel mundial el ictus es la segunda causa de mortalidad y la tercera más común en los países industrializados [2].
- La mortalidad hospitalaria se sitúa en torno al 12.9 % [3].

¿Por qué es un problema?

- El ictus isquémico supone la segunda causa de mortalidad en nuestro país en la población general [1].
- La primera causa de mortalidad en la mujer [1].
- A nivel mundial el ictus es la segunda causa de mortalidad y la tercera más común en los países industrializados [2].
- La mortalidad hospitalaria se sitúa en torno al 12.9 % [3].

¿Por qué es un problema?

- El ictus isquémico supone la segunda causa de mortalidad en nuestro país en la población general [1].
- La primera causa de mortalidad en la mujer [1].
- A nivel mundial el ictus es la segunda causa de mortalidad y la tercera más común en los países industrializados [2].
- La mortalidad hospitalaria se sitúa en torno al 12.9 % [3].

¿Por qué es un problema?

- El ictus isquémico supone la segunda causa de mortalidad en nuestro país en la población general [1].
- La primera causa de mortalidad en la mujer [1].
- A nivel mundial el ictus es la segunda causa de mortalidad y la tercera más común en los países industrializados [2].
- La mortalidad hospitalaria se sitúa en torno al 12.9 % [3].

¿Por qué se puede/debe estudiar?

- Se han detectado una serie de factores de riesgo que permiten estimar la probabilidad de fallecer o presentar secuelas [3].
- Ya existen trabajos que permiten desarrollar modelos predictivos para estimar dicha mortalidad [4][5].
- Un plan integral de actuaciones desde la llegada del paciente aumenta las probabilidades de recuperación de este.
- El tiempo, las medidas, la formación y el equipo que presta atención se ha demostrado ser un factor determinante en la tasa de mortalidad intrahospitalaria.

¿Por qué se puede/debe estudiar?

- Se han detectado una serie de factores de riesgo que permiten estimar la probabilidad de fallecer o presentar secuelas [3].
- Ya existen trabajos que permiten desarrollar modelos predictivos para estimar dicha mortalidad [4][5].
- Un plan integral de actuaciones desde la llegada del paciente aumenta las probabilidades de recuperación de este.
- El tiempo, las medidas, la formación y el equipo que presta atención se ha demostrado ser un factor determinante en la tasa de mortalidad intrahospitalaria.

¿Por qué se puede/debe estudiar?

- Se han detectado una serie de factores de riesgo que permiten estimar la probabilidad de fallecer o presentar secuelas [3].
- Ya existen trabajos que permiten desarrollar modelos predictivos para estimar dicha mortalidad [4][5].
- Un plan integral de actuaciones desde la llegada del paciente aumenta las probabilidades de recuperación de este.
- El tiempo, las medidas, la formación y el equipo que presta atención se ha demostrado ser un factor determinante en la tasa de mortalidad intrahospitalaria.

¿Por qué se puede/debe estudiar?

- Se han detectado una serie de factores de riesgo que permiten estimar la probabilidad de fallecer o presentar secuelas [3].
- Ya existen trabajos que permiten desarrollar modelos predictivos para estimar dicha mortalidad [4][5].
- Un plan integral de actuaciones desde la llegada del paciente aumenta las probabilidades de recuperación de este.
- El tiempo, las medidas, la formación y el equipo que presta atención se ha demostrado ser un factor determinante en la tasa de mortalidad intrahospitalaria.

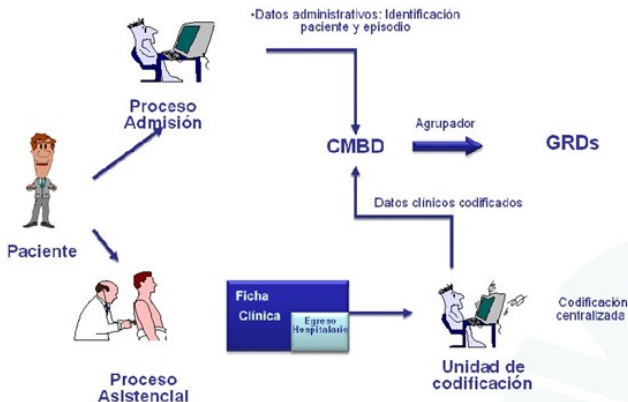
Hipótesis del trabajo

- Es posible identificar determinadas variables contenidas en el CMBD asociadas a la mortalidad intrahospitalaria de pacientes con ictus no lisado.
- Es posible elaborar un modelo predictivo de mortalidad intrahospitalaria basado en dichas variables.

Hipótesis del trabajo

- Es posible identificar determinadas variables contenidas en el CMBD asociadas a la mortalidad intrahospitalaria de pacientes con ictus no lisado.
- Es posible elaborar un modelo predictivo de mortalidad intrahospitalaria basado en dichas variables.

CMBD



Objetivos del proyecto y TFM

- Validación externa, Recalibración de los coeficientes, Trasladar y transferir a todos los niveles asistenciales implicados en la atención del ictus este modelo.

TFM

- Realizar una revisión metodológica de los pasos seguidos. Respalando los resultados con una metodología diferente.
- Traer nuevas técnicas utilizadas en la Ciencia de Datos.
- Desarrollo de nuevos modelos de estimación de mortalidad con los que se espera igualar o mejorar los resultados, tanto para **pre-ingreso** del paciente, como para **post-ingreso**, lo que se corresponderá con los **escenarios A y B**.

Objetivos del proyecto y TFM

- Validación externa, Recalibración de los coeficientes, Trasladar y transferir a todos los niveles asistenciales implicados en la atención del ictus este modelo.

TFM

- Realizar una revisión metodológica de los pasos seguidos. Respaldando los resultados con una metodología diferente.
- Traer nuevas técnicas utilizadas en la Ciencia de Datos.
- Desarrollo de nuevos modelos de estimación de mortalidad con los que se espera igualar o mejorar los resultados, tanto para **pre-ingreso** del paciente, como para **post-ingreso**, lo que se corresponderá con los **escenarios A y B**.

Objetivos del proyecto y TFM

- Validación externa, Recalibración de los coeficientes, Trasladar y transferir a todos los niveles asistenciales implicados en la atención del ictus este modelo.

TFM

- Realizar una revisión metodológica de los pasos seguidos. Respaldando los resultados con una metodología diferente.
- Traer nuevas técnicas utilizadas en la Ciencia de Datos.
- Desarrollo de nuevos modelos de estimación de mortalidad con los que se espera igualar o mejorar los resultados, tanto para **pre-ingreso** del paciente, como para **post-ingreso**, lo que se corresponderá con los **escenarios A y B**.

Objetivos del proyecto y TFM

- Validación externa, Recalibración de los coeficientes, Trasladar y transferir a todos los niveles asistenciales implicados en la atención del ictus este modelo.

TFM

- Realizar una revisión metodológica de los pasos seguidos. Respalando los resultados con una metodología diferente.
- Traer nuevas técnicas utilizadas en la Ciencia de Datos.
- Desarrollo de nuevos modelos de estimación de mortalidad con los que se espera igualar o mejorar los resultados, tanto para **pre-ingreso** del paciente, como para **post-ingreso**, lo que se corresponderá con los **escenarios A y B**.

Índice

1 Introducción

- ¿Por qué estudiar este problema?
- CMBD: Conjunto Mínimo Básico de Datos
- Objetivos del proyecto

2 Análisis Exploratorio de Datos

- 2 escenarios
- Recorrido por las variables más importantes
- Prevalencia de los predictores
- Discusión de la métrica y división muestra

3 Metodología

- Datos desbalanceados
- Algoritmo Genético de Selección de Predictores

4 Escenarios

- Escenario A
 - Modelos Utilizados
- Escenario B
 - Modelos Utilizados
- Conclusiones y Trabajo Futuro

*"La principal diferencia entre un Data Scientist **Junior** y un **Senior** radica en que, mientras uno prueba distintos modelos para cercar el problema, el otro se comprueba qué está ocurriendo en los datos."*

2 escenarios

- **Escenario A:** Se realizará todo el proceso utilizando únicamente variables seleccionadas por criterio experto.
- **Escenario B:** Nos preguntamos si podemos mejorar los resultados obtenidos en A, utilizando toda la base de datos, y algoritmos de selección de variables de ciencia de datos.

Pasos a seguir

- Normalizar nombres de columnas.
- Eliminación de observaciones con algún valor faltante.
- Transformación de variables categóricas y booleanas.

2 escenarios

- **Escenario A:** Se realizará todo el proceso utilizando únicamente variables seleccionadas por criterio experto.
- **Escenario B:** Nos preguntamos si podemos mejorar los resultados obtenidos en A, utilizando toda la base de datos, y algoritmos de selección de variables de ciencia de datos.

Pasos a seguir

- Normalizar nombres de columnas.
- Eliminación de observaciones con algún valor faltante.
- Transformación de variables categóricas y booleanas.

2 escenarios

- **Escenario A:** Se realizará todo el proceso utilizando únicamente variables seleccionadas por criterio experto.
- **Escenario B:** Nos preguntamos si podemos mejorar los resultados obtenidos en A, utilizando toda la base de datos, y algoritmos de selección de variables de ciencia de datos.

Pasos a seguir

- Normalizar nombres de columnas.
- Eliminación de observaciones con algún valor faltante.
- Transformación de variables categóricas y booleanas.

2 escenarios

- **Escenario A:** Se realizará todo el proceso utilizando únicamente variables seleccionadas por criterio experto.
- **Escenario B:** Nos preguntamos si podemos mejorar los resultados obtenidos en A, utilizando toda la base de datos, y algoritmos de selección de variables de ciencia de datos.

Pasos a seguir

- Normalizar nombres de columnas.
- Eliminación de observaciones con algún valor faltante.
- Transformación de variables categóricas y booleanas.

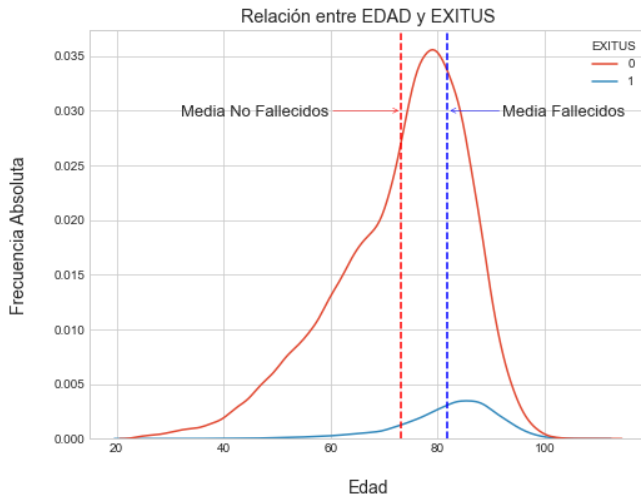
2 escenarios

- **Escenario A:** Se realizará todo el proceso utilizando únicamente variables seleccionadas por criterio experto.
- **Escenario B:** Nos preguntamos si podemos mejorar los resultados obtenidos en A, utilizando toda la base de datos, y algoritmos de selección de variables de ciencia de datos.

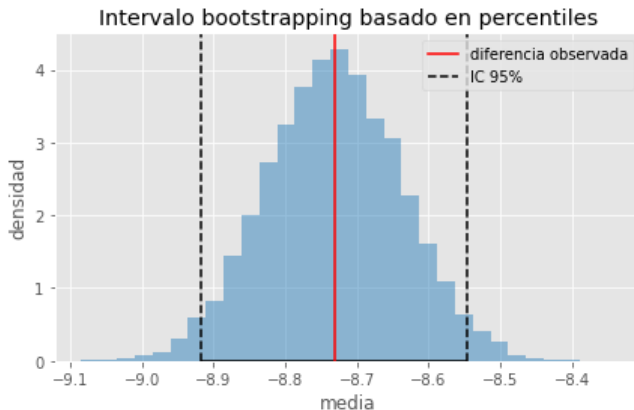
Pasos a seguir

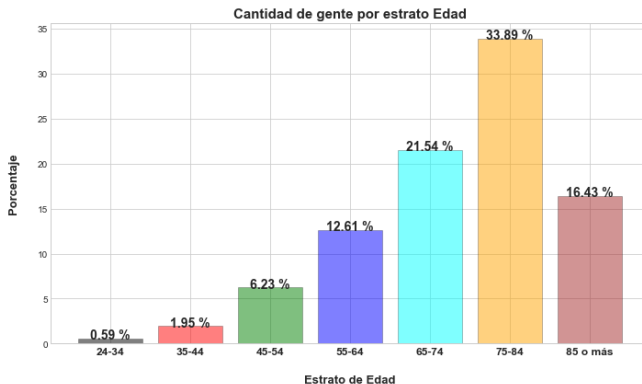
- Normalizar nombres de columnas.
- Eliminación de observaciones con algún valor faltante.
- Transformación de variables categóricas y booleanas.

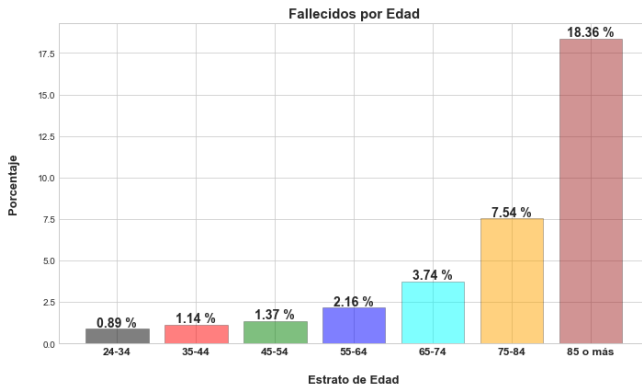
Edad



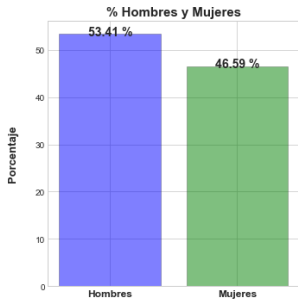
Diferencia edad entre grupos



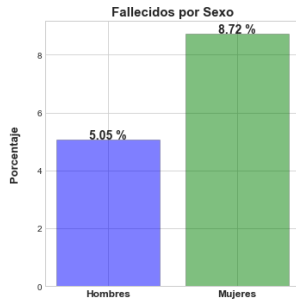




Fallecidos por sexo

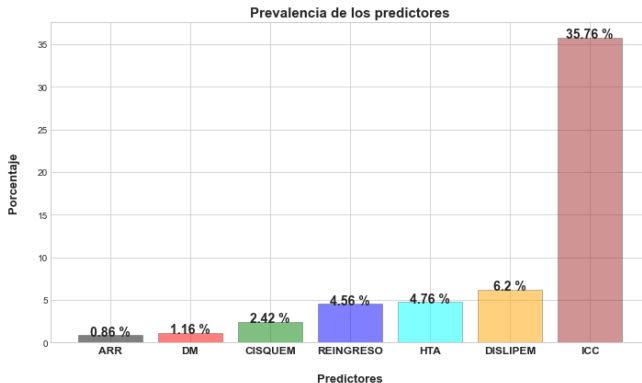


(k)

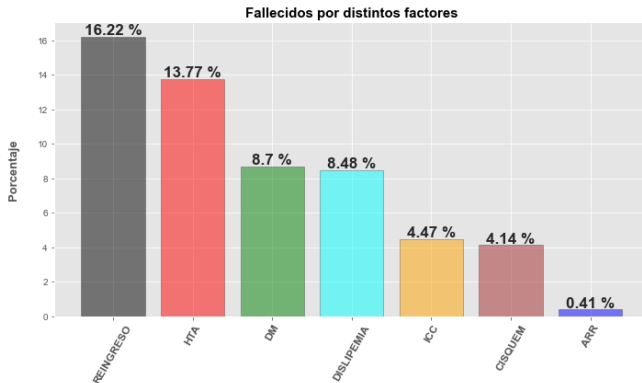


(l)

Prevalencia de los predictores



Fallecidos por predictor



- Debido a la distribución de la variable objetivo (93 % - 7 %), no tiene sentido usar la métrica **Accuracy**.
- Necesitamos minimizar la cantidad de falsos negativos, lo que equivale a optimizar el **recall**.
- También es importante no descuidar la precisión (optimización costo-eficiencia). Pero es mucho peor tener falsos positivos que falsos negativos.
- La curva **ROC** nos ayudará a estudiar la habilidad de discriminación del modelo entre las 2 clases.
- Mediremos también el **F1-Score** y veremos cómo se comporta con la matriz de confusión.
- Incluso **balanced-acc**:

$$\frac{TNR + TPR}{2}$$

- Debido a la distribución de la variable objetivo (93 % - 7 %), no tiene sentido usar la métrica **Accuracy**.
- Necesitamos minimizar la cantidad de falsos negativos, lo que equivale a optimizar el **recall**.
- También es importante no descuidar la precisión (optimización costo-eficiencia). Pero es mucho peor tener falsos positivos que falsos negativos.
- La curva **ROC** nos ayudará a estudiar la habilidad de discriminación del modelo entre las 2 clases.
- Mediremos también el **F1-Score** y veremos cómo se comporta con la matriz de confusión.
- Incluso **balanced-acc**:

$$\frac{TNR + TPR}{2}$$

- Debido a la distribución de la variable objetivo (93 % - 7 %), no tiene sentido usar la métrica **Accuracy**.
- Necesitamos minimizar la cantidad de falsos negativos, lo que equivale a optimizar el **recall**.
- También es importante no descuidar la precisión (optimización costo-eficiencia). Pero es mucho peor tener falsos positivos que falsos negativos.
- La curva **ROC** nos ayudará a estudiar la habilidad de discriminación del modelo entre las 2 clases.
- Mediremos también el **F1-Score** y veremos cómo se comporta con la matriz de confusión.
- Incluso **balanced-acc**:

$$\frac{TNR + TPR}{2}$$

- Debido a la distribución de la variable objetivo (93 % - 7 %), no tiene sentido usar la métrica **Accuracy**.
- Necesitamos minimizar la cantidad de falsos negativos, lo que equivale a optimizar el **recall**.
- También es importante no descuidar la precisión (optimización costo-eficiencia). Pero es mucho peor tener falsos positivos que falsos negativos.
- La curva **ROC** nos ayudará a estudiar la habilidad de discriminación del modelo entre las 2 clases.
- Mediremos también el **F1-Score** y veremos cómo se comporta con la matriz de confusión.
- Incluso **balanced-acc**:

$$\frac{TNR + TPR}{2}$$

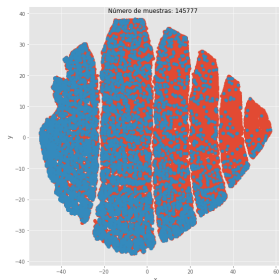
- Debido a la distribución de la variable objetivo (93 % - 7 %), no tiene sentido usar la métrica **Accuracy**.
- Necesitamos minimizar la cantidad de falsos negativos, lo que equivale a optimizar el **recall**.
- También es importante no descuidar la precisión (optimización costo-eficiencia). Pero es mucho peor tener falsos positivos que falsos negativos.
- La curva **ROC** nos ayudará a estudiar la habilidad de discriminación del modelo entre las 2 clases.
- Mediremos también el **F1-Score** y veremos cómo se comporta con la matriz de confusión.
- Incluso **balanced-acc**:

$$\frac{TNR + TPR}{2}$$

- Debido a la distribución de la variable objetivo (93 % - 7 %), no tiene sentido usar la métrica **Accuracy**.
- Necesitamos minimizar la cantidad de falsos negativos, lo que equivale a optimizar el **recall**.
- También es importante no descuidar la precisión (optimización costo-eficiencia). Pero es mucho peor tener falsos positivos que falsos negativos.
- La curva **ROC** nos ayudará a estudiar la habilidad de discriminación del modelo entre las 2 clases.
- Mediremos también el **F1-Score** y veremos cómo se comporta con la matriz de confusión.
- Incluso **balanced-acc**:

$$\frac{TNR + TPR}{2}$$

- Finalmente, tenemos 36 variables, división train 85 % - test 15 % (establecemos `random_state = 44`, estratificada).
- Método de entrenamiento y evaluación: **Validación cruzada con Holdout** (GridSearchCV).



(m) escenario A



(n) escenario B

Índice

1 Introducción

- ¿Por qué estudiar este problema?
- CMBD: Conjunto Mínimo Básico de Datos
- Objetivos del proyecto

2 Análisis Exploratorio de Datos

- 2 escenarios
- Recorrido por las variables más importantes
- Prevalencia de los predictores
- Discusión de la métrica y división muestra

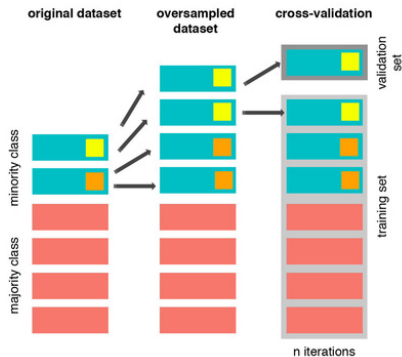
3 Metodología

- Datos desbalanceados
- Algoritmo Genético de Selección de Predictores

4 Escenarios

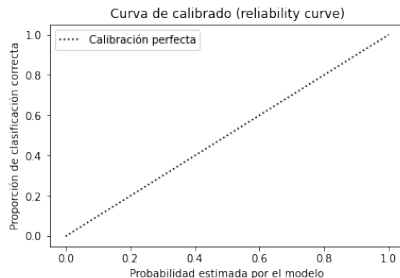
- Escenario A
 - Modelos Utilizados
- Escenario B
 - Modelos Utilizados
- Conclusiones y Trabajo Futuro

- Ignorar el problema.
- Submuestrear la clase mayoritaria.
- Sobremuestrear la clase minoritaria.



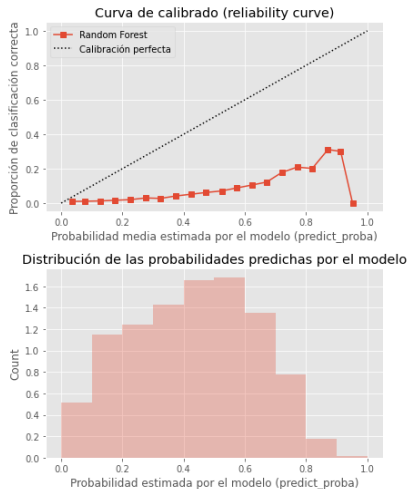
Calibrado del modelo

Un modelo calibrado es aquel en el que, el valor estimado de probabilidad, puede interpretarse directamente como la confianza que se tiene de que la clasificación predicha es correcta.



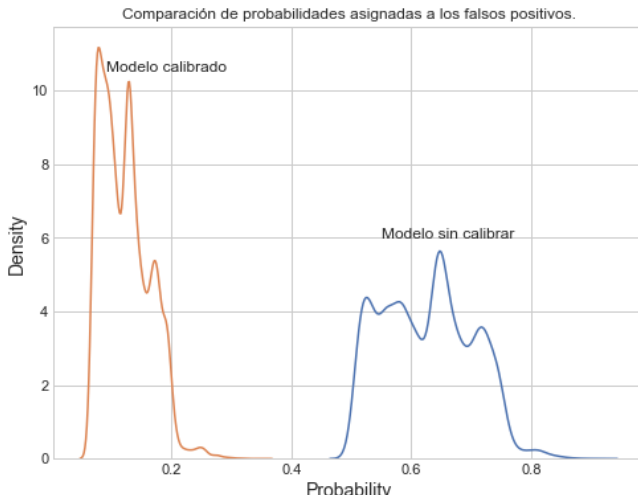
Objetivo

- Hacer buenas estimaciones del riesgo de mortalidad.



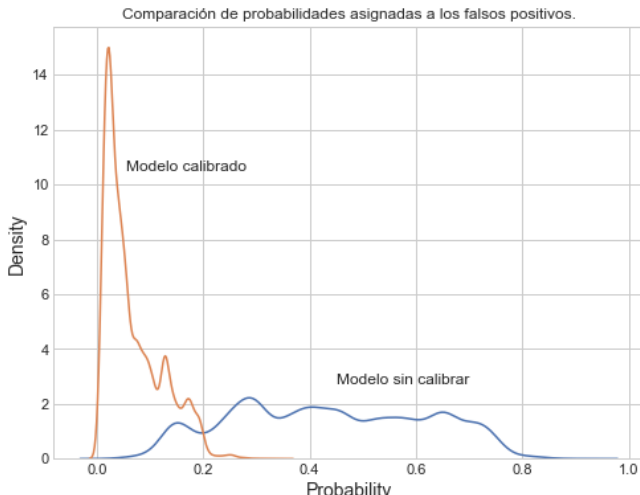
Modelo calibrado vs normal

- Estimaciones sobre falsos positivos.

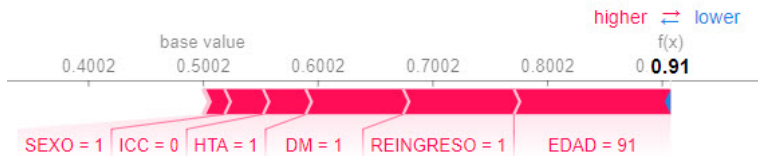


Modelo calibrado vs normal

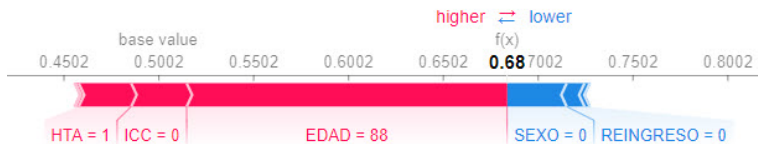
- Distribución de probabilidad total.



¿Falsos positivos?



(ñ)



(o)

- **Objetivo:** Encontrar la combinación de predictores que da lugar al mejor modelo (randomforest, f1).
- **Funcionamiento:**
 - Creación de población de P individuos (combinación de predictores).
 - Calcular fitness de cada individuo (métrica de calidad).
 - Crear una población vacía y repetir los pasos hasta crear P nuevos individuos:
 - Seleccionar dos individuos de la población existente
 - Cruzar los dos individuos seleccionados para generar un nuevo descendiente (crossover).
 - Aplicar un proceso de mutación aleatorio sobre el nuevo individuo.
 - Añadir el nuevo individuo a la nueva población.
 - Reemplazar la antigua población de la nueva.
 - Si no se cumple el criterio de parada, volver al paso 2.

- **Objetivo:** Encontrar la combinación de predictores que da lugar al mejor modelo (randomforest, f1).
- **Funcionamiento:**
 - Creación de población de P individuos (combinación de predictores).
 - Calcular fitness de cada individuo (métrica de calidad).
 - Crear una población vacía y repetir los pasos hasta crear P nuevos individuos:
 - Seleccionar dos individuos de la población existente
 - Cruzar los dos individuos seleccionados para generar un nuevo descendiente (crossover).
 - Aplicar un proceso de mutación aleatorio sobre el nuevo individuo.
 - Añadir el nuevo individuo a la nueva población.
 - Reemplazar la antigua población de la nueva.
 - Si no se cumple el criterio de parada, volver al paso 2.

- **Objetivo:** Encontrar la combinación de predictores que da lugar al mejor modelo (randomforest, f1).
- **Funcionamiento:**
 - Creación de población de P individuos (combinación de predictores).
 - Calcular fitness de cada individuo (métrica de calidad).
 - Crear una población vacía y repetir los pasos hasta crear P nuevos individuos:
 - Seleccionar dos individuos de la población existente
 - Cruzar los dos individuos seleccionados para generar un nuevo descendiente (crossover).
 - Aplicar un proceso de mutación aleatorio sobre el nuevo individuo.
 - Añadir el nuevo individuo a la nueva población.
 - Reemplazar la antigua población de la nueva.
 - Si no se cumple el criterio de parada, volver al paso 2.

- **Objetivo:** Encontrar la combinación de predictores que da lugar al mejor modelo (randomforest, f1).
- **Funcionamiento:**
 - Creación de población de P individuos (combinación de predictores).
 - Calcular fitness de cada individuo (métrica de calidad).
 - Crear una población vacía y repetir los pasos hasta crear P nuevos individuos:
 - Seleccionar dos individuos de la población existente
 - Cruzar los dos individuos seleccionados para generar un nuevo descendiente (crossover).
 - Aplicar un proceso de mutación aleatorio sobre el nuevo individuo.
 - Añadir el nuevo individuo a la nueva población.
 - Reemplazar la antigua población de la nueva.
 - Si no se cumple el criterio de parada, volver al paso 2.

- **Objetivo:** Encontrar la combinación de predictores que da lugar al mejor modelo (randomforest, f1).
- **Funcionamiento:**
 - Creación de población de P individuos (combinación de predictores).
 - Calcular fitness de cada individuo (métrica de calidad).
 - Crear una población vacía y repetir los pasos hasta crear P nuevos individuos:
 - Seleccionar dos individuos de la población existente
 - Cruzar los dos individuos seleccionados para generar un nuevo descendiente (crossover).
 - Aplicar un proceso de mutación aleatorio sobre el nuevo individuo.
 - Añadir el nuevo individuo a la nueva población.
 - Reemplazar la antigua población de la nueva.
 - Si no se cumple el criterio de parada, volver al paso 2.

Índice

1 Introducción

- ¿Por qué estudiar este problema?
- CMBD: Conjunto Mínimo Básico de Datos
- Objetivos del proyecto

2 Análisis Exploratorio de Datos

- 2 escenarios
- Recorrido por las variables más importantes
- Prevalencia de los predictores
- Discusión de la métrica y división muestra

3 Metodología

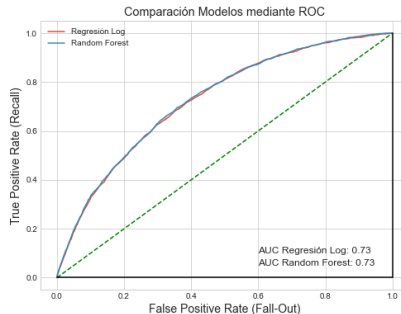
- Datos desbalanceados
- Algoritmo Genético de Selección de Predictores

4 Escenarios

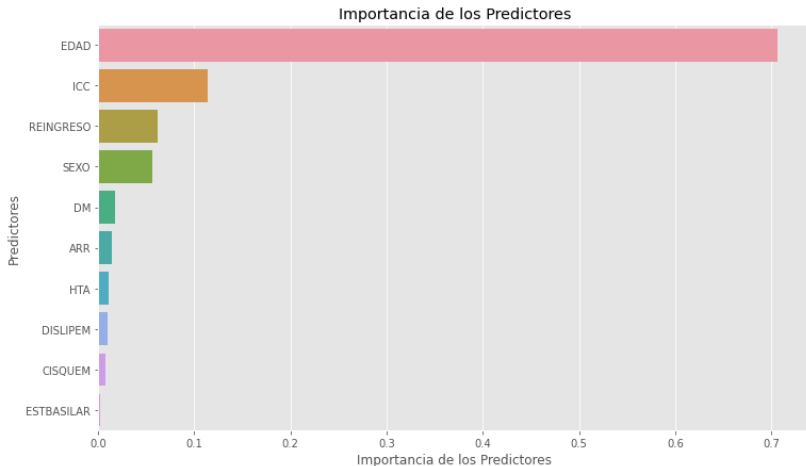
- Escenario A
 - Modelos Utilizados
- Escenario B
 - Modelos Utilizados
- Conclusiones y Trabajo Futuro

Modelos Utilizados

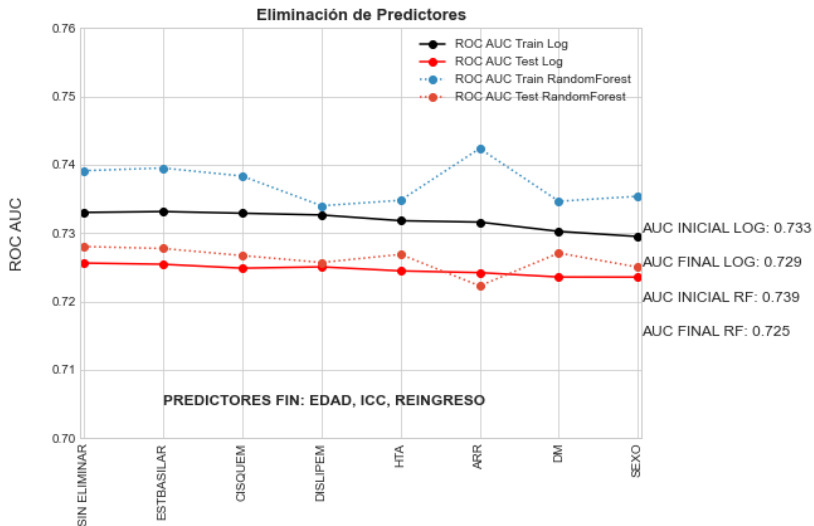
- Regresión logística balanceada.
- Pipeline Random Undersample + Regresión logística.
- Random Forest Balanceado.
- Pipeline Random Undersample + Random Forest.



Importancia de los predictores



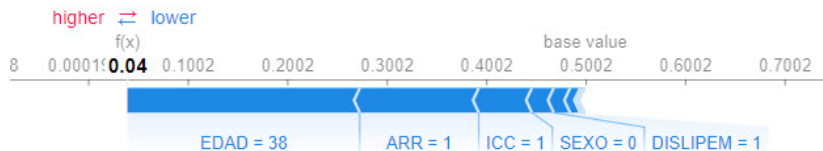
Eliminación Recursiva de predictores



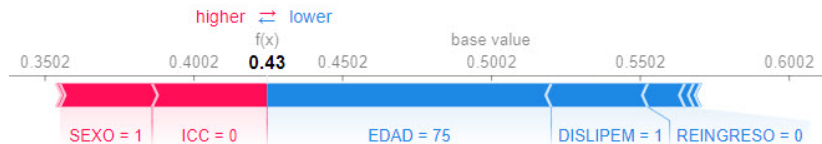
Métricas de los modelos

	Accuracy	Balanced accuracy	F1-Score	AUROC	Precision	Recall
Balanced Logistic Regression	0.633756	0.671532	0.208851	0.732837	0.122282	0.715214
Random Undersample + Logistic Regression	0.637426	0.671994	0.209760	0.732947	0.123003	0.711966
Balanced random forest	0.651886	0.673818	0.213537	0.733811	0.126013	0.699178
Under-sampling + Random forest	0.639648	0.661089	0.204715	0.712233	0.120337	0.685882

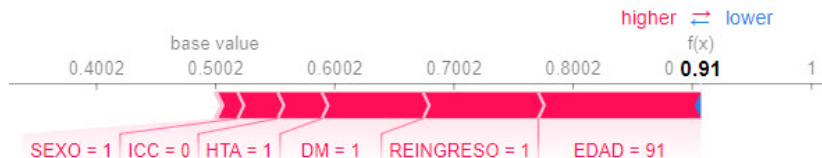
Métricas de los modelos



Métricas de los modelos

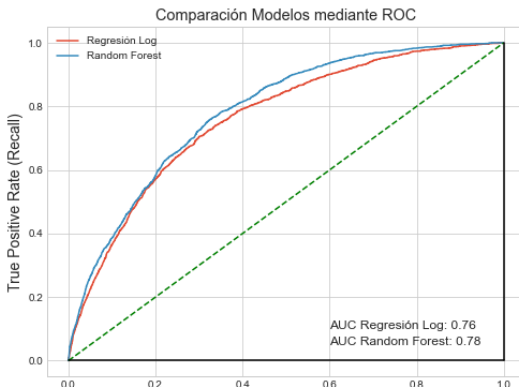


Métricas de los modelos



Variables del algoritmo selección genética

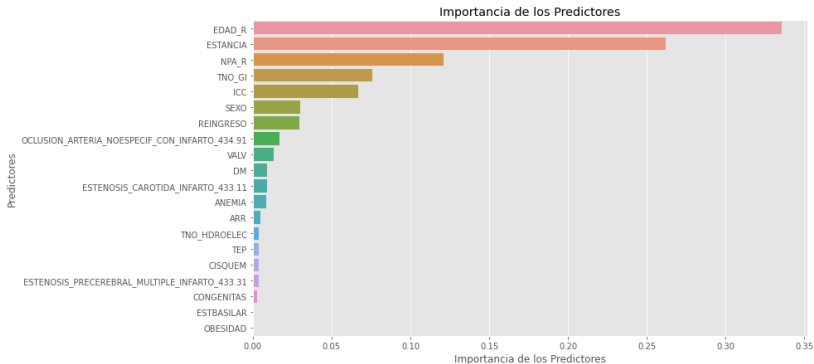
- Regresión logística balanceada.
- Pipeline Random Undersample + Regresión logística.
- Random Forest Balanceado.
- Pipeline Random Undersample + Random Forest.



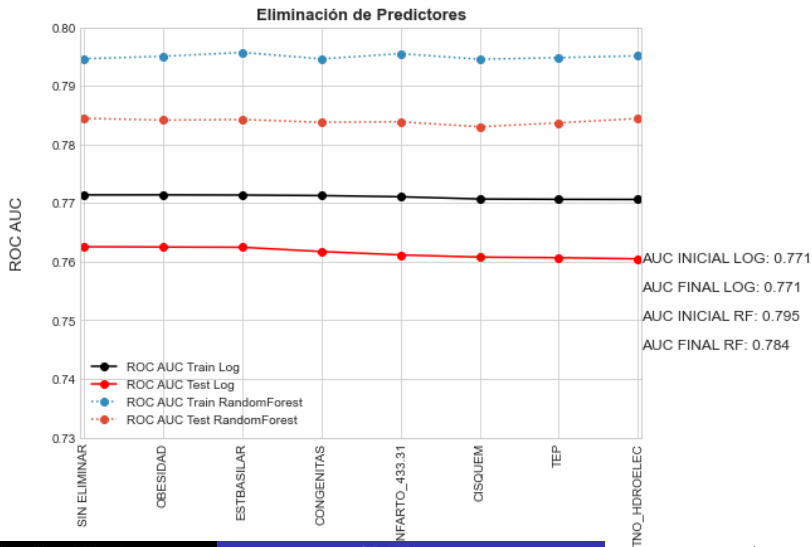
Resultado modelos

	Accuracy	Bal-Acc	F1-Score	AUROC	Precis.	Recall
Balanced Logistic Regression	0.680313		0.235936	0.770848		0.730236
		0.703465			0.140699	
Random Undersample + Logistic Regression	0.683414		0.236999	0.770645		0.727394
		0.703810			0.141567	
Balanced random forest	0.688579		0.245163	0.788986		0.748200
		0.716228			0.146609	
Under-sampling + Random forest	0.671052		0.220933	0.737112		0.690045
		0.679860			0.131527	

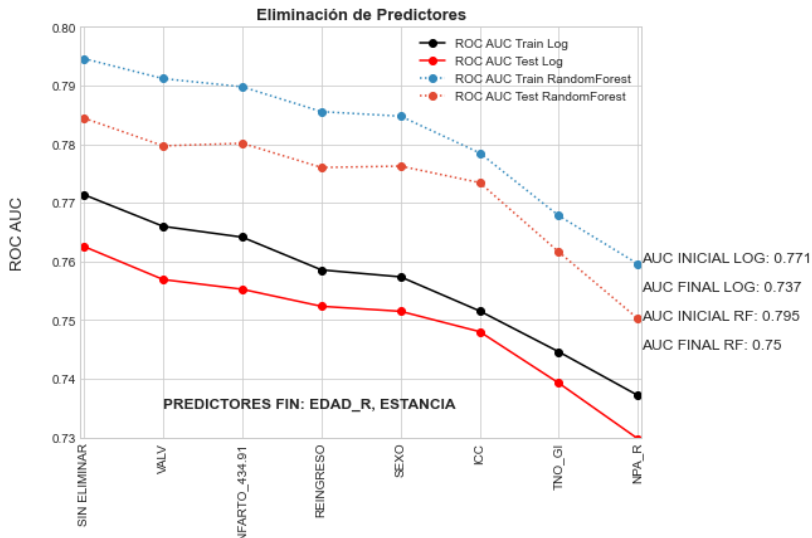
Importancia modelo



Eliminación recursiva



Eliminación recursiva



Conclusiones

- Hemos conseguido desarrollar un modelo que teniendo el mismo número de variables que el del escenario A, mejora los resultados.
- Eliminar variables cuya importancia del modelo es menor a 0,01 no afecta al resultado de éste.
- ¿Qué modelo poner en producción?
- Servicios AWS
 - SageMaker.
 - API Gateway, S3, Lambdas, DynamoDB.

Conclusiones

- Hemos conseguido desarrollar un modelo que teniendo el mismo número de variables que el del escenario A, mejora los resultados.
 - Eliminar variables cuya importancia del modelo es menor a 0,01 no afecta al resultado de éste.
-
- ¿Qué modelo poner en producción?
 - Servicios AWS
 - SageMaker.
 - API Gateway, S3, Lambdas, DynamoDB.

Conclusiones

- Hemos conseguido desarrollar un modelo que teniendo el mismo número de variables que el del escenario A, mejora los resultados.
 - Eliminar variables cuya importancia del modelo es menor a 0,01 no afecta al resultado de éste.
-
- ¿Qué modelo poner en producción?
 - Servicios AWS
 - SageMaker.
 - API Gateway, S3, Lambdas, DynamoDB.

Conclusiones

- Hemos acercado técnicas propias de Ciencia de Datos a la medicina, donde se espera que tengan una buena acogida y una futura aplicación.
- Se ha realizado una aproximación a un problema que está abierto.
- De cara al futuro se espera poder desplegar todo lo desarrollado de manera que se puedan seguir optimizando estos modelos. Usando nuevos datos con los que validaremos y recalibraremos los modelos.

Bibliografía I

Nos hemos basado en los siguientes documentos:



Joaquín Amat.

Algoritmo genético de selección de predictores.

Web,

https://www.cienciadedatos.net/documentos/py03_eleccion_predictores



Comes E Oliveres M Targa C Balcells M et al Arboix A,
García-Eroles L.

*Importancia del perfil cardiovascular en la mortalidad hospitalaria
de los infartos cerebrales.*

Revista española de cardiología, 2008.

Bibliografía II



Ministerio de Sanidad.

Estrategia en Ictus del Sistema Nacional de Salud: Ministerio de Sanidad y Política Social.

Editorial Universidad de Almería, España, 2009.



Bennett DA Anderson CS. Feigin VL, Lawes CM.

"Stroke epidemiology: a review of population-based studies of incidence, prevalence, and case-fatality in the late 20th century".

The Lancet Neurology, 2003.



Whisnant JP.

Modeling of risk factors for ischemic stroke.

The Willis Lecture. Stroke; a journal of cerebral circulation., 2013.

Bibliografía III



Scikit-learn.

Documentación de Scikit-learn.

Scikit-learn, <https://scikit-learn.org/stable/>, 2021.



Dai D Olson DM Reeves MJ Saver JL et al Smith EE, Shobha N.
A risk score for in-hospital death in patients admitted with ischemic or hemorrhagic stroke.

Journal of the American Heart Association, 2013.



Juan Manuel García Torrecillas.

Memoria de Registro del Modelo a Validar, Modelo predictivo para la mortalidad hospitalaria en el ictus isquémico no lisado.

Hospital Torrecárdenas, 2020.

¡Muchas gracias por su
atención!