

Modelo predictivo para validación de datos sintéticos clínicos en diabetes

Javier I. Ramírez^{1,6}, Francisco M. Ortuño^{1,2}, Carlos Loucera^{2,3}, Joaquin Dopazo^{2,3,4,5}

¹ Departamento de Ingeniería de Computadores, Automática y Robótica, Universidad de Granada (UGR), Granada, España.

² Plataforma de Medicina Computacional, Fundación Progreso y Salud (FPS), Sevilla, España.

³ Medicina Computacional de Sistemas, Instituto de Biomedicina de Sevilla (IBiS), Sevilla, España.

⁴ Bioinformática en Enfermedades Raras (BiER), Centro de Investigación Biomédica en Red de Enfermedades Raras (CIBERER), Seville, Spain.

⁵ FPS/ELIXIR-ES, Virgen del Rocío Hospital, Seville, Spain.

⁶ FIDESOL- Fundación en Investigación y Desarrollo del Software Libre, Granada, España.

OBJETIVO

En este trabajo se desarrolla un sistema de predicción de patologías médicas utilizando datos secuenciales de diagnósticos de pacientes diabéticos. La meta principal es mejorar la precisión en la detección temprana de diversas enfermedades, lo que podría tener un impacto significativo en la toma de decisiones clínicas y en la prevención de complicaciones graves. Este estudio se centra en el análisis de datos médicos temporales generados por una GAN basada en datos reales obtenidos por la BPS, permitiendo capturar las dinámicas y patrones que preceden a la aparición de patologías.

Dataset pacientes diabéticos

	Cantidad
Pacientes	1.215.974
Códigos Patologías	83

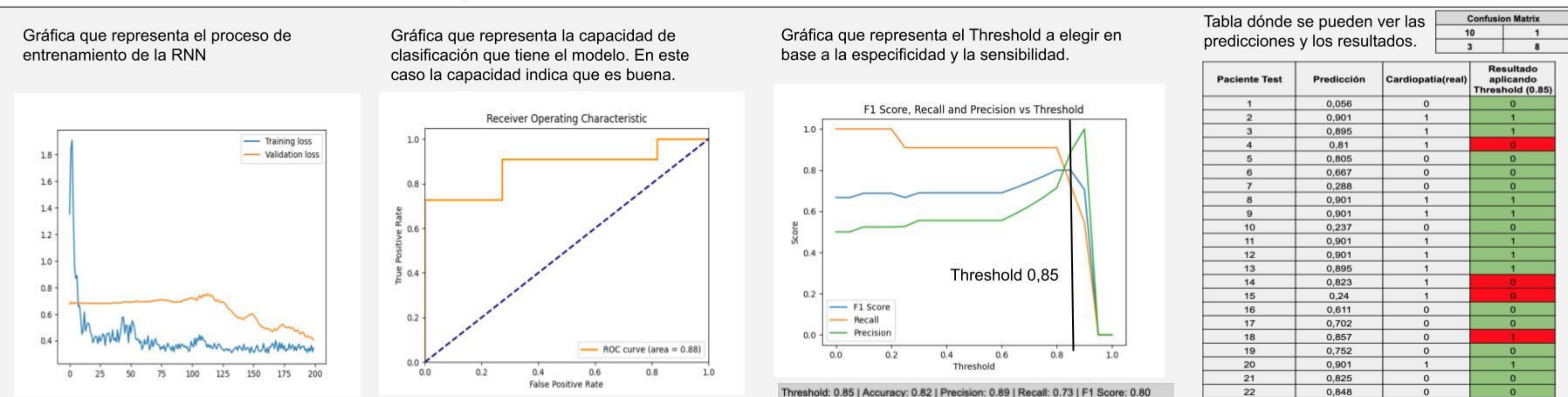
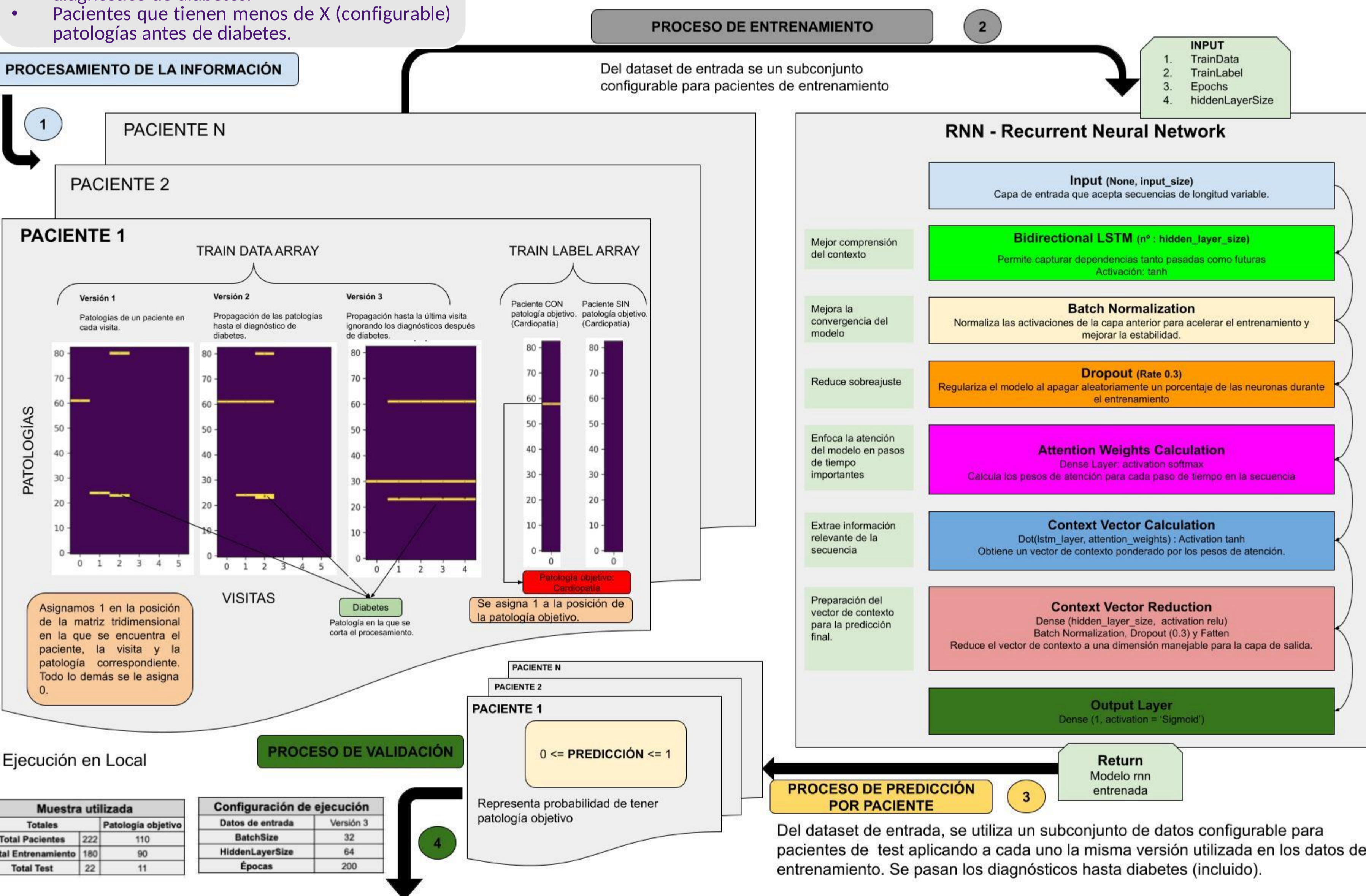
Dataset proporcionado por uno de los retos propuestos por el congreso anual CAMDA (Critical Assessment of Massive Data Analysis). En este reto, se proporciona un dataset sintético que trazan las trayectorias de diagnóstico de pacientes diabéticos. Estos registros fueron creados mediante Autoencoders Duales Adversarios entrenados con datos de 1,2 millones de pacientes reales de la Base Poblacional de Salud (BPS) del Sistema Andaluz de Salud.



Contiene información ordenada en el tiempo sobre las enfermedades crónicas diagnosticadas en cada una de las visitas que haya realizado cada paciente. Además, también hay información sobre las décadas en las que se han hecho cada visita.

- PACIENTES ELIMINADOS:**
- Pacientes sin diabetes
 - Pacientes con la patología objetivo antes del diagnóstico de diabetes.
 - Pacientes que tienen menos de X (configurable) patologías antes de diabetes.

METODOLOGÍA



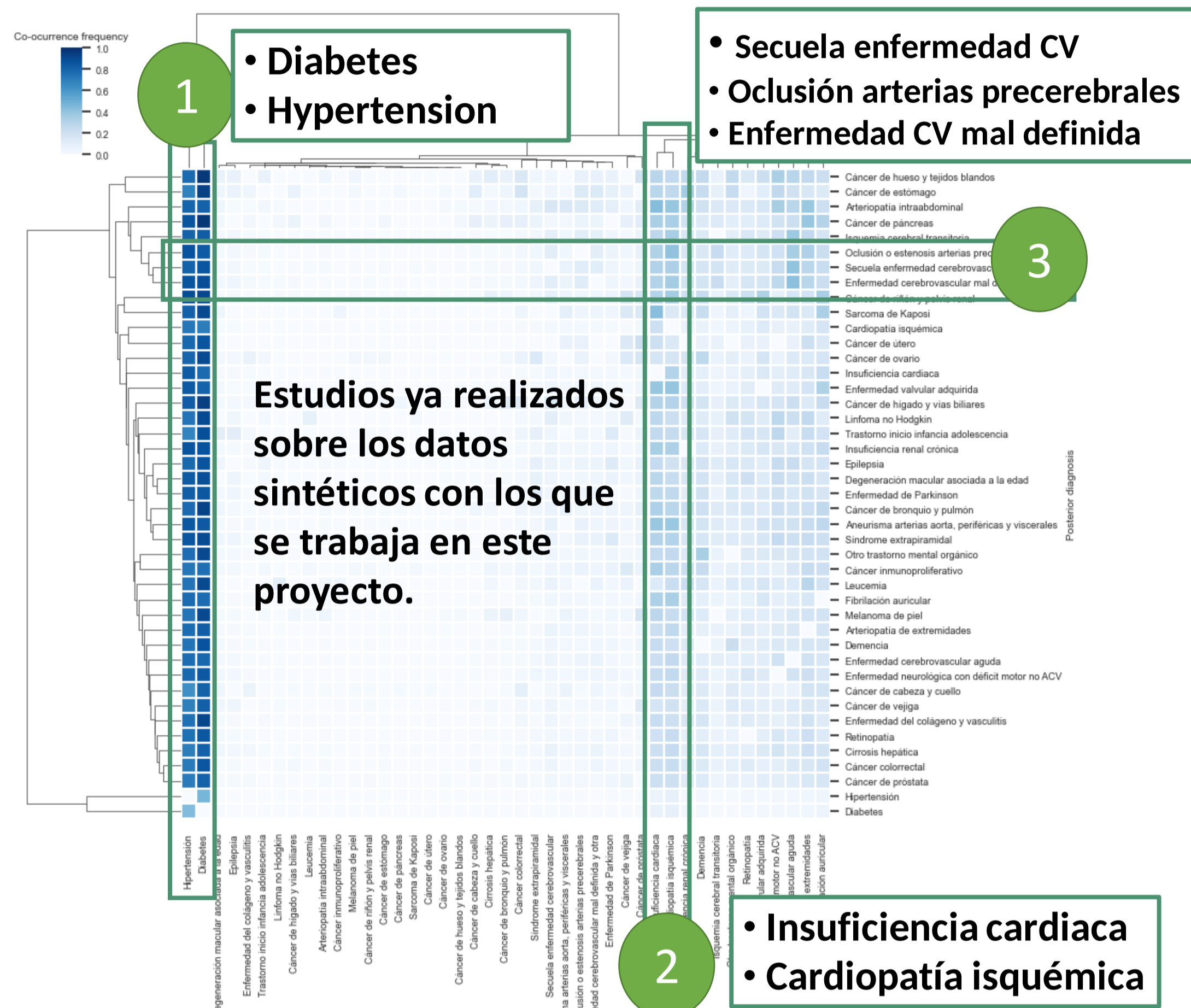
Pruebas funcionales Dataset completo (Cluster)

Ejecución	Versión	Numero Patologías previas a Diabetes	Total Pacientes	Pacientes Entrenamiento	Pacientes Test	Balanceada patología objetivo	Batchsize	Épocas	Hidden Layer Size	ROC	F1-Score
1	1	1	65718	52576	13142	si	64	600	64	0,55	0,60
2	2	3	318807	255045	63762	no	64	600	64	0,55	0,10
3	2	5	7724	6950	774	si	128	600	128	0,58	0,67
4	2	3	62254	56028	6226	si	128	600	128	0,56	0,67
5	3	1	177692	159922	17770	si	128	600	128	0,59	0,67
6	3	2	111162	100044	11118	si	128	600	128	0,61	0,67
7	3	3	62254	56028	6226	si	128	600	128	0,62	0,67
8	3	4	32420	29178	3242	si	128	600	128	0,6	0,67

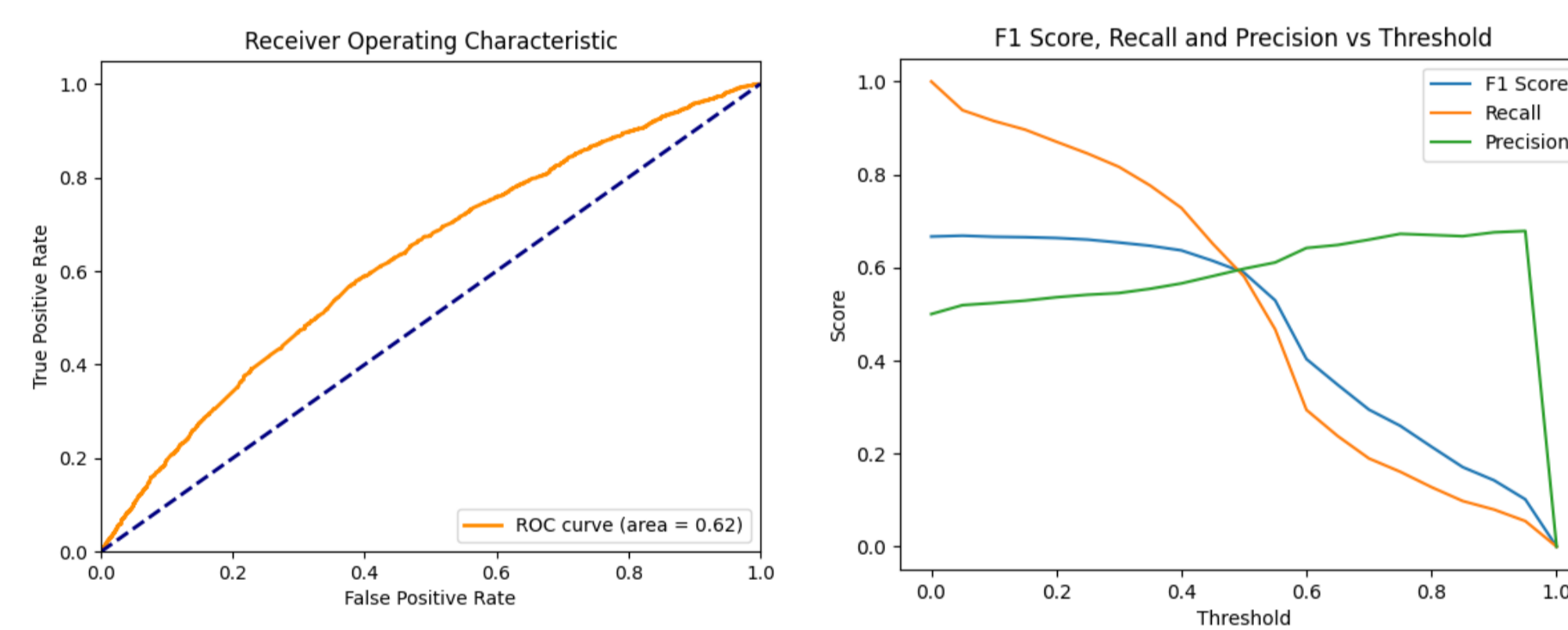
Los resultados muestran que todavía hace falta ajustar el modelo para el dataset completo. De todas formas, para este tipo de problemas parece que la RNN se comporta bien para muestras en las que hay un número significativo de patologías antes de Diabetes. Por otra parte, puede que la patología objetivo elegida no tenga ningún patrón dentro del dataset con lo que habría que probar con otras patologías objetivo.

CONCLUSION

Secuencias de co-ocurrencias más frecuentes en dataset



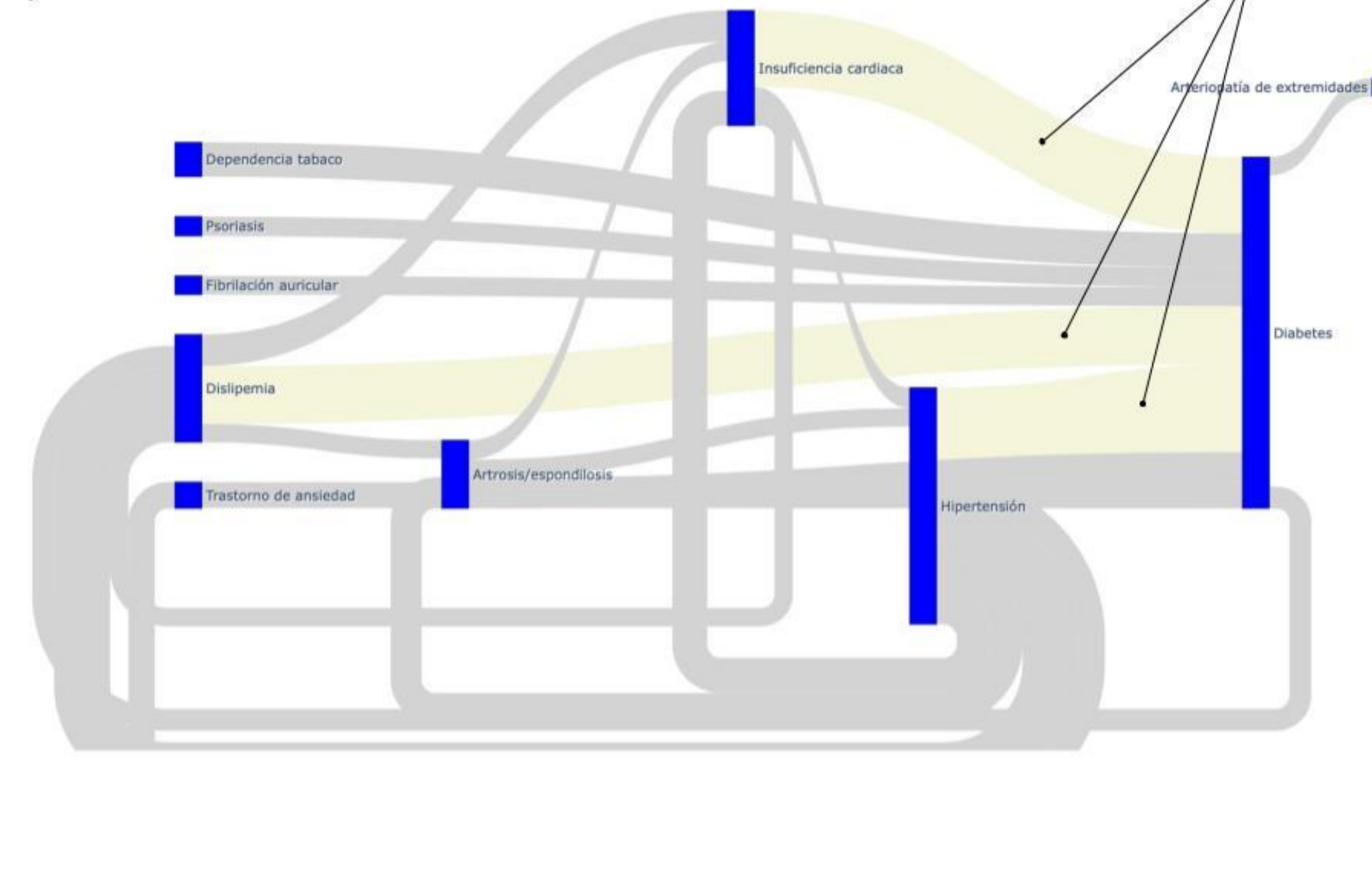
Mejor Resultado en Cluster



Transiciones Frecuentes - TOP 10
Se utilizan pacientes que tienen más de 0.8 de predicción.

Pacientes que según la predicción tienen más de un 80% de tener Cardiopatía

Top 3 transiciones más frecuentes



Transiciones Frecuentes - TOP 10
Se utilizan pacientes que tienen menos de 0.3 de predicción.

Pacientes que según la predicción tienen menos de un 30% de tener Cardiopatía

