

La elección metodológica entre enfoques tradicionales y aprendizaje automático en salud pública

The Methodological Choice Between Traditional Approaches and Machine Learning in Public Health

Maicel Monzón Pérez^{1*}

[0000-0003-2117-9145](#)

¹Escuela Nacional de Salud Pública. Cuba.

* Autor para la correspondencia: maicel@infomed.sld.cu

RESUMEN

En el campo de la bioestadística, existe un debate metodológico en torno a la distinción entre dos propósitos principales del modelado: modelos explicativos y modelos predictivos.

En el presente trabajo se exploran brevemente los dos paradigmas: uno centrado en la generación de conocimiento y el otro en el desarrollo y validación de tecnologías para la toma de decisiones informadas y se realiza una aproximación a la modelación predictiva en la clínica mediante el aprendizaje automático. Se incluyen ejemplos en código R, útiles para la implementación de problemas similares.

Palabras clave: inferencia causal; modelos diagnósticos; modelos pronósticos; tecnologías para la toma de decisiones informadas.

ABSTRACT

In the field of biostatistics, there is a methodological debate surrounding the distinction between two main purposes of modeling: explanatory models and predictive models.

This paper briefly explores the two fundamental paradigms underlying modeling in public health: one focused on knowledge generation and the other on the development and validation of technologies for informed decision-making. It also provides an approach to predictive modeling in clinical practice using machine learning. Examples in R code, which are useful for implementing similar problems, are also included.

Keywords: causal inference; diagnostic models; prognostic models; technologies for informed decision-making.

Recibido: 13/08/2025

Aprobado: 15/12/2025



Introducción

En el campo de la bioestadística existe un debate metodológico fundamental que gira en torno a la distinción entre dos propósitos principales del modelado: los modelos diseñados para explicar fenómenos en una población, y aquellos orientados a predecir características o eventos futuros en un individuo. Estos últimos se conocen comúnmente como modelos de predicción clínica. ⁽¹⁾

Silva Ayçaguer (2012) plantea que los modelos explicativos priorizan la inferencia causal y la interpretabilidad de los resultados, lo que se alinea con el propósito general de la investigación científica de comprender mecanismos subyacentes; Steyerberg (2009) plantea que los modelos predictivos están más vinculados al desarrollo de tecnologías. ^{(1),(2)} Por ejemplo, su objetivo es construir herramientas fiables para estimar la probabilidad de que ciertas enfermedades estén presentes (modelos diagnósticos) o predecir resultados futuros (modelos pronósticos), lo cual es crucial para la toma de decisiones clínicas.

Los modelos de predicción clínica se pueden enmarcar dentro del enfoque de aprendizaje automático (machine learning), dado que comparten objetivos, métricas de evaluación y principios metodológicos fundamentales. ⁽³⁾ No obstante, en la práctica aplicada resulta frecuente la confusión entre ambos paradigmas, especialmente cuando se emplean algoritmos similares con finalidades conceptualmente distintas (Fig. 1).

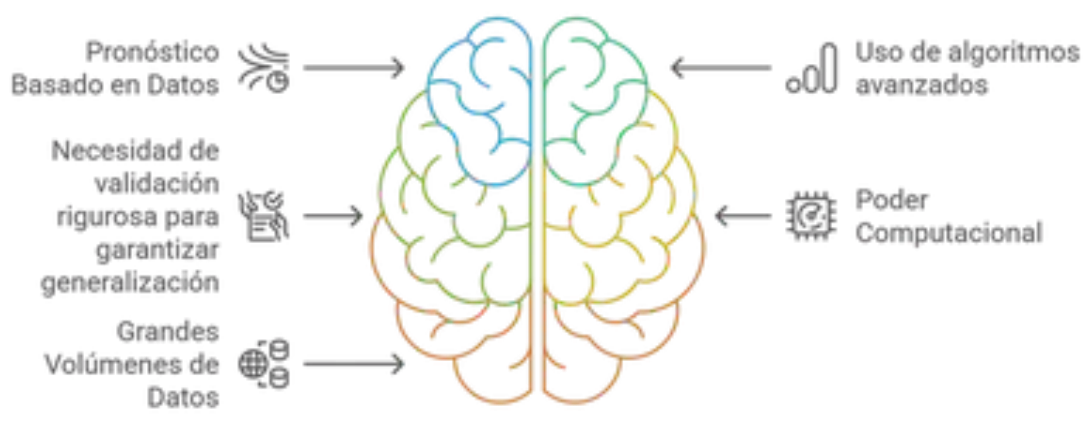


Fig.1- Elementos clave del aprendizaje automático aplicados a la predicción clínica.

El presente trabajo no pretende realizar una comparación técnica exhaustiva entre la bioestadística clásica y el aprendizaje automático aplicado a la salud. Su objetivo consiste en explorar de forma concisa los dos paradigmas fundamentales que subyacen al modelado en salud pública -inferencia y predicción- y en señalar las consecuencias metodológicas y prácticas derivadas de su confusión. Para ello se presentan ejemplos ilustrativos, mediante código en R, útiles para la implementación problemas análogos.



Desarrollo

Antes de implementar cualquier modelo resulta imprescindible definir con claridad el objetivo fundamental del análisis, ya que esta decisión condiciona la metodología empleada, el uso de los datos y los criterios de evaluación adoptados. ⁽¹⁾

El enfoque explicativo se orienta a comprender relaciones entre variables y a cuantificar asociaciones en una población específica. El modelo funciona como una herramienta para generar hipótesis y contribuir a la explicación de mecanismos, y la pregunta central se formula en términos de magnitud e incertidumbre de la asociación entre un predictor y un desenlace.

El enfoque predictivo, en cambio, se centra en la construcción de herramientas capaces de generalizar adecuadamente a datos nuevos y no observados. En este caso el modelo se concibe como un medio para apoyar decisiones informadas sobre individuos en contextos distintos al original. La pregunta clave se orienta a estimar la probabilidad de ocurrencia del desenlace y a evaluar la fiabilidad de dicha estimación.

Existe la percepción de que el aprendizaje automático se limita a algoritmos complejos de tipo caja negra, como los Bosques Aleatorios o las Redes Neuronales Profundas. Sin embargo, este campo incluye también modelos altamente interpretables, como la regresión logística, que pueden explicarse tanto con fines explicativos como predictivos, en función del marco metodológico adoptado.

Para ilustrar ambos enfoques se utiliza el conjunto de datos "Pima indians diabetes", incluido en el paquete mlbench, ampliamente utilizado con fines metodológicos y educativos. ⁽⁵⁾ La preparación de los datos incorpora la imputación múltiple como parte integral del proceso analítico, en reconocimiento de que la omisión sistemática de valores faltantes introduce sesgos que afectan tanto la inferencia como la predicción.

Bloque 1: Preparación del entorno e imputación

```
library(tidyverse)
```

```
library(tidymodels)
```

```
library(missRanger)
```

```
library(mlbench)
```

```
library(rms)
```

```
data(PimaIndiansDiabetes2)
```

Imputación múltiple para evitar sesgos de selección

```
datos <- PimaIndiansDiabetes2 %>%
```

```
  missRanger(num.trees = 100, verbose = 0) %>%
```

```
  mutate(diabetes = if_else(diabetes == "pos", 1, 0)) %>%
```

```
  select(diabetes, age, pedigree, pregnant, mass, glucose, pressure)
```



Enfoque explicativo

Cuando el objetivo del análisis se orienta a la explicación, el foco se sitúa en la estimación de efectos y en la cuantificación de la incertidumbre asociada. En este contexto se emplea el conjunto completo de datos, con el propósito de maximizar la precisión de las estimaciones.⁽²⁾

Bloque 3: Validación

```
set.seed(1353)
```

```
split <- initial_split(datos, strata = "diabetes")
```

```
entrenamiento <- training(split)
```

```
prueba <- testing(split)
```

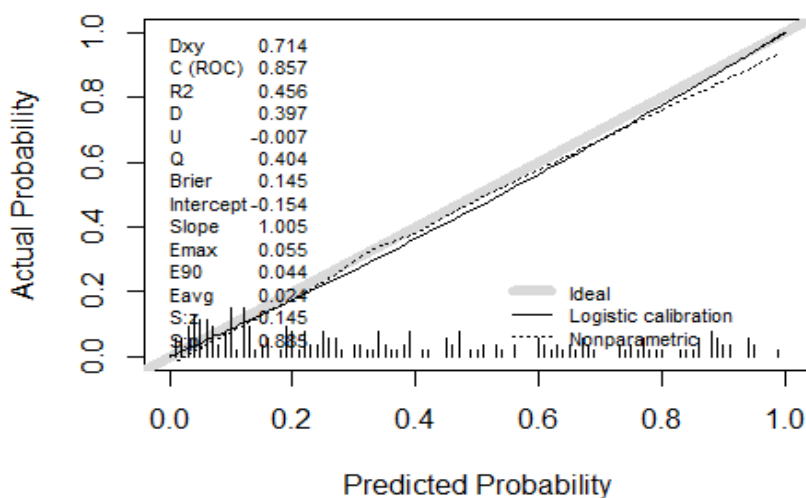
Ajuste en entrenamiento y evaluación en prueba

```
modelo_predictivo <- lrm(diabetes ~ ., data = entrenamiento, x = TRUE, y = TRUE)
```

```
predicciones <- predict(modelo_predictivo, newdata = prueba, type = "fitted")
```

Calibración y Discriminación

```
val.prob(predicciones, prueba$diabetes)
```



##	Dxy	C (ROC)	R2	D	D:Chi-sq	D:p
##	0.713671642	0.856835821	0.456410383	0.397099673	77.243137157	0.000000000
##	U	U:Chi-sq	U:p	Q	Brier	Intercept
##	-0.006905713	0.674103147	0.713872022	0.404005385	0.145260812	-0.154105017
##	Slope	Emax	E90	Eavg	S:z	S:p
##	1.004694638	0.054661953	0.044242321	0.024006879	-0.145030589	0.884686720



Los coeficientes obtenidos mediante regresión logística cuantifican asociaciones condicionadas a los predictores incluidos y a los supuestos del modelo. Por ejemplo, el índice de predisposición genética (*pedigree*) muestra una asociación estadísticamente significativa con la presencia de diabetes, expresada mediante un Odds Ratio superior a la unidad.

En estudios observacionales, estas asociaciones no permiten establecer relaciones causales directas. La posible presencia de confusión residual, sesgo de selección y variables no medidas limita la interpretación causal de los resultados (1,6). En consecuencia, los hallazgos adquieren valor como evidencia asociativa y como base para la generación de hipótesis, pero no como demostración de causalidad. En este enfoque, el modelo no se concibe como una herramienta diagnóstica, sino como un instrumento para la comprensión estadística del fenómeno estudiado.

Enfoque predictivo

En el modelado predictivo, la interpretación aislada de los coeficientes pierde relevancia frente a la evaluación del rendimiento global del modelo en datos no utilizados durante el entrenamiento. La separación estricta entre conjuntos de entrenamiento y prueba responde a la necesidad de estimar de forma honesta la capacidad de generalización.

La evaluación del desempeño predictivo requiere considerar tanto la discriminación como la calibración. Un valor elevado del estadístico C indica una buena capacidad para diferenciar entre individuos con y sin el desenlace; sin embargo, una alta discriminación no garantiza decisiones clínicas seguras si las probabilidades estimadas no reflejan riesgos reales. La calibración adquiere así un papel central, dado que una desalineación sistemática entre riesgos predichos y observados puede inducir errores en la práctica clínica.

La inclusión de variables estrechamente relacionadas con la definición del desenlace, como los niveles de glucosa en el caso de la diabetes, introduce una forma de circularidad diagnóstica. Esta situación puede generar una sobreestimación del rendimiento predictivo y una falsa percepción de eficacia. En este contexto, dicha inclusión cumple un propósito ilustrativo, al mostrar la necesidad de interpretar con cautela las métricas de desempeño cuando se emplean predictores de alto peso diagnóstico.

Los modelos predictivos no constituyen entidades estáticas. Cambios en la prevalencia de la enfermedad, en las características de la población o en las prácticas clínicas pueden deteriorar su rendimiento con el tiempo. Por ello, la validación debe concebirse como un proceso continuo que incluye la monitorización sistemática y, cuando resulte necesario, la recalibración periódica del modelo.

A continuación, en la tabla 1 se resumen las diferencias entre ambos modelos.



Tabla 1- Diferenciación entre modelos explicativos y predictivos.

Dimensión	Modelo Explicativo (Inferencia)	Modelo Predictivo (Tecnología)
Pregunta Clave	¿Cuál es la relación entre X e Y?	¿Cuál es el riesgo real de Y para este individuo?
Uso de Datos	Dataset completo para estabilidad de la estimación.	Separación estricta (Entrenamiento / Prueba).
Salida Principal	Coefficientes (OR) e Intervalos de Confianza.	Métricas de rendimiento (AUC, Calibración).
Peligro Crítico	Sobreinterpretación causal de datos observacionales.	Sobreajuste (Overfitting) y falsa seguridad clínica.

Conclusiones

La bioestadística moderna y la ciencia de datos responsable constituyen enfoques complementarios dentro del conocimiento basado en datos. La discusión metodológica relevante no se sitúa en la comparación entre algoritmos específicos, sino en la distinción entre inferencia y predicción. Para el investigador clínico y el epidemiólogo, el rigor inferencial continúa siendo esencial para la generación de conocimiento y la formulación de hipótesis. Para el clínico o el gestor de salud que requiere herramientas de apoyo a la decisión, la validación predictiva rigurosa resulta ineludible. Confundir inferencia con predicción no representa un error técnico menor, sino una falla conceptual que puede conducir tanto a conclusiones científicas erróneas como a decisiones clínicas inseguras.

Referencias

1. Steyerberg EW. Clinical Prediction Models. Philadelphia: Springer; 2009.

2. Ayçaguer LC, Suárez P. Cultura estadística e investigación científica en el campo de la salud. España: Ediciones Díaz de Santos; 1998.

3. Ryu L, Han K. Machine learning vs. statistical models for prediction modelling. J Korean Soc Radiol. 2022;83(6):1219.

4. Monzón Pérez M. Cómo entrenar y validar un modelo de machine learning [Internet]. España: Bioestadística edu; 2025 [citado 2025 Dic 27]. Disponible en: <https://maicel.netlify.app/post/como-entrenar-y-validar-un-modelo-de-machine-learning/>

5. Leisch F, Dimitriadou E. mlbench: Machine Learning Benchmark Problems. En: Harrell FE. Regression Modeling Strategies. Philadelphia: Springer; 2015.

Conflicto de interés

No se declaran conflictos de interés.

Declaración de autoría

La responsabilidad de la concepción, redacción y revisión del manuscrito corresponde íntegramente al autor.

