



Evaluación de la eficiencia con modelos no paramétricos

PRÁCTICA 1

Análisis de eficiencia y Productividad

Enunciado

Para la realización de la práctica utilizaremos la base de datos “rice_producers.xlsx” disponible en el [CAMPUS VIRTUAL UMH](#). Además, también se utilizarán los siguientes 3 *toy datasets* (A, B y C, respectivamente):

Store	A	B	C	D	E	F	G	H
x_1 : employee	2	3	3	4	5	5	6	8
y_1 : sale	1	3	2	3	4	2	3	5

Store	A	B	C	D	E	F	G	H	I
x_1 : employee	4	7	8	4	2	5	6	5.5	6
x_2 : floor area	3	3	1	2	4	2	4	2.5	2.5
y_1 : sale	1	1	1	1	1	1	1	1	1

Store	A	B	C	D	E	F	G
x_1 : employee	1	1	1	1	1	1	1
y_1 : customers	1	2	3	4	4	5	6
y_2 : sales	5	7	4	3	6	5	2

También se ha simulado una base de datos (D) con 50 DMUs generada como se muestra a continuación:

$$x_1 \sim U(a=1, b=10)$$

$$u \sim |N(\mu=1, \sigma=0.4)|$$

$$y_D = \ln(x_1) + 3$$

$$y_1 = y_D - u$$

Preguntas teóricas

En el contexto del **Análisis de Eficiencia** se desea evaluar la **eficiencia** de una muestra de n unidades llamadas *Decision-Making Units* [DMUs], donde cada DMU_i , $i=1,\dots,n$, consume $\mathbf{x}_i = (x_{i1}, \dots, x_{ij}, \dots, x_{im}) \in \mathbb{R}_+^m$ **inputs** para la producción de $\mathbf{y}_i = (y_{i1}, \dots, y_{ir}, \dots, y_{is}) \in \mathbb{R}_+^s$ **outputs**.

Se asume que las DMUs son generadas a partir de un Proceso Generador de Datos (PGD). En el caso de considerar un único *output*, el PGD es una función **desconocida**, **monótona no decreciente** y generalmente **cóncava**:

$$f(\mathbf{x}): \mathbb{R}_+^m \rightarrow \mathbb{R}_+$$

Este PGD $[f(\mathbf{x})]$ se conoce como la **frontera teórica de producción** y mide cuál es el máximo output producible dando cierto nivel de recursos. Por ejemplo, ¿cuál es el máximo número de zapatos (y_1) que se pueden fabricar dado cierto número de trabajadores ($x_1 = 5$)?

La estimación de esta frontera de producción (llamada **frontera de Mejores Prácticas** en algunos contextos) y la medición de la eficiencia de las unidades de la muestra puede llevarse a cabo bajo dos metodologías bien diferenciadas: **enfoques paramétricos** y **enfoques no paramétricos**.

- Un modelo es considerado **paramétrico** cuando el número de parámetros a estimar es **fijo** y determinado a priori.
- Un modelo es considerado **no paramétrico** cuando el número de parámetros a estimar **no es fijo** y viene determinado por la muestra de datos, los hiperparámetros que definen el modelo, etc.

La principal diferencia entre ambas metodologías es la **presunción previa** de una forma funcional del PGD. Por ejemplo, bajo un enfoque paramétrico, podemos considerar que $f(\mathbf{x})$ es una función de producción de tipo Cobb-Douglas (generalmente, es su forma log-lineal):

$$y_D = \alpha_0 \cdot x_1^{\alpha_1} \cdot x_2^{\alpha_2} \cdot x_3^{\alpha_3}$$

$$\ln(y_D) = \ln(\alpha_0 \cdot x_1^{\alpha_1} \cdot x_2^{\alpha_2} \cdot x_3^{\alpha_3})$$

$$\ln(y_D) = \ln(\alpha_0) + \alpha_1 \cdot \ln(x_1) + \alpha_2 \cdot \ln(x_2) + \alpha_3 \cdot \ln(x_3),$$

y estimar, entonces, el vector de coeficientes $\alpha = (\alpha_0, \alpha_1, \alpha_2, \alpha_3)$ que hace que la expresión $\ln(\alpha_0) + \alpha_1 \cdot \ln(x_1) + \alpha_2 \cdot \ln(x_2) + \alpha_3 \cdot \ln(x_3)$ se ajuste lo mejor posible a los datos disponibles (al *output* observado). En el enfoque paramétrico, el número de parámetros a estimar suele ser relativamente pequeño en relación al número de *inputs*. Además, este tipo de enfoques suelen permitir interpretaciones claras de los modelos, por ejemplo, en el caso lineal, α_1 representa el cambio marginal del máximo *output* que se puede producir si modificamos el primer *input* dejando constante el resto de los *inputs*. Si estamos ante un modelo **log-log** (como el de arriba) aumentar $x_1 \rightarrow e^1 \cdot x_1$ producirá un aumento en el valor esperado de $y \rightarrow e^{\alpha_1} \cdot y$. En el enfoque no paramétrico, a priori, no se asume ninguna forma funcional concreta para $f(\mathbf{x})$.

Finalmente, en cuanto a los *outputs* observados (y_i), cabe resaltar que son traslaciones (verticales) de este PGD:

$$y_i = f(\mathbf{x}_i) - u_i + \varepsilon_i, \quad i = 1, \dots, n$$

donde u mide la ineficiencia técnica de una DMU y ε mide cierto error aleatorio.

- Los **modelos paramétricos** generalmente asumen que $u \sim |N(0, \sigma_u)|$ y $\varepsilon \sim N(0, \sigma_\varepsilon)$. Debido a que partimos de la existencia de error aleatorio, se les denomina modelos **estocásticos**. Estos supuestos permiten estimar los parámetros del modelo mediante métodos como el de máxima verosimilitud y realizar inferencia estadística sobre los mismos: intervalos de confianza y/o contrastes de hipótesis sobre su significatividad.
- Los **modelos no paramétricos** únicamente asumen que $u_i \geq 0, i = 1, \dots, n$ y no consideran error aleatorio, es decir, $\varepsilon_i = 0, i = 1, \dots, n$. Precisamente, debido a que no se considera la existencia de error aleatorio, se dice que este tipo de modelos son **deterministas**. Por lo tanto, bajo este enfoque, cierta *DMU_i* será técnicamente eficiente cuando $u_i = 0$.

Además, en el caso del Análisis Envolvente de Datos (modelo no paramétrico en el que nos centraremos), siempre hablaremos de **eficiencia técnica “relativa”**, dado que dependerá exclusivamente de la muestra de datos utilizada.

En la asignatura de Análisis de Eficiencia y Productividad nos centramos en el estudio y aplicación de modelos no paramétricos, desde su enfoque tradicional hasta los recientes avances desde el campo del Aprendizaje Automático.

Ejercicio 1. Dado el siguiente conjunto de datos:

Empresa	Inversión en I+D	Empleados	Ingresos Totales	Patentes
Innovatech S.A.	5.0	105	100	4
Soluciones Digitales S.L.	2.5	125	60	7
TecnoDynamics Inc.	8.0	50	150	5
CiberInnovación Ltda.	3.0	80	80	6
InfoTech Global	10.0	160	200	3
Onda Digital	1.5	150	50	5
Sistemas Inteligentes	4.0	140	110	4
AlphaData Tech	6.5	90	130	6

- ¿Qué variables son de tipo *input* y cuáles son de tipo *output*? Justifica tu respuesta.
- Dadas las variables “inversión en I+D” e “ingresos totales”, identifica, mediante una representación gráfica, que unidades de la muestra son técnicamente eficientes.
- Dadas las variables “inversión en I+D” y “empleados”, identifica, mediante una representación gráfica, que unidades de la muestra son técnicamente eficientes.
- Dadas las variables “ingresos totales” y “patentes”, identifica, mediante una representación gráfica, que unidades de la muestra son técnicamente eficientes.
- Proporciona una expresión que permita calcular la eficiencia de las unidades

Ejercicio 2. Menciona que ventajas y/o inconvenientes presentan cada una de las metodologías frente a la otra.

Ejercicio 3. Consideramos una muestra de n DMUs, para las cuales se desea evaluar la eficiencia técnica. Cada DMU consume $\mathbf{x}_i = (x_{i1}, \dots, x_{ij}, \dots, x_{im}) \in \mathbb{R}_+^m$ *inputs* para producir $\mathbf{y}_i = (y_{i1}, \dots, y_{ir}, \dots, y_{is}) \in \mathbb{R}_+^s$ *outputs*. Para medir la eficiencia (relativa) de cada DMU, es necesario definir un conjunto tecnológico común T compartido por todas las DMUs de la muestra. Desde una perspectiva más amplia, esta tecnología puede expresarse como:

$$T = \left\{ (\mathbf{x}, \mathbf{y}) \in \mathbb{R}_+^{m+s} : \mathbf{x} \text{ puede producir } \mathbf{y} \right\}.$$

Existen tres tipos de tecnologías de producción:

Nombre	Dataset	Convexidad	Libre disponibilidad	Rendimientos
CCR	<i>Toy dataset A</i>	✓	✓	Constantes
BCC	<i>Toy dataset B</i>	✓	✓	Variables
FDH	<i>Toy dataset C</i>	✗	✓	Variables

A continuación, se define formalmente cada una de estas tecnologías. A partir de su definición, proporciona un punto que pertenezca a cada una de las tecnologías para cada una de las bases de datos que se proporcionan en la tabla anterior.

Tecnología DEA (CCR)

La tecnología CCR (Charnes, A., Cooper, W. W., & Rhodes, E. (1978). Measuring the efficiency of decision making units. *European journal of operational research*, 2(6), 429-444.) es una tecnología convexa que satisface el principio de libre disponibilidad bajo rendimientos constantes a escala:

$$\hat{T}_{CRS} = \left\{ (\mathbf{x}, \mathbf{y}) \in \mathbb{R}_+^{m+s} : x_j \geq \sum_{i=1}^n \lambda_i x_{ij}, j=1, \dots, m, y_r \leq \sum_{i=1}^n \lambda_i y_{ir}, r=1, \dots, s, \lambda_i \geq 0, i=1, \dots, n \right\}.$$

Tecnología DEA (BCC)

La tecnología BCC (Banker, R. D., Charnes, A., & Cooper, W. W. (1984). Some Models for Estimating Technical and Scale Inefficiencies in Data Envelopment Analysis. *Management Science*, 30(9), 1078-1092.) es una tecnología convexa que satisface el principio de libre disponibilidad bajo rendimientos variables a escala:

$$\hat{T}_{VRS} = \left\{ (\mathbf{x}, \mathbf{y}) \in \mathbb{R}_+^{m+s} : x_j \geq \sum_{i=1}^n \lambda_i x_{ij}, j=1, \dots, m, y_r \leq \sum_{i=1}^n \lambda_i y_{ir}, r=1, \dots, s, \sum_{i=1}^n \lambda_i = 1, \lambda_i \geq 0, i=1, \dots, n \right\}.$$

Tecnología FDH

La tecnología FDH (Deprins, D., Simar, L., & Tulkens, H. (2006). Measuring labor-efficiency in post offices. In *Public goods, environmental externalities and fiscal competition* (pp. 285-

309). Boston, MA: Springer US.) es una tecnología no convexa que únicamente satisface el principio de libre disponibilidad:

$$\hat{T}_{FDH} = \left\{ (\mathbf{x}, \mathbf{y}) \in \mathbb{R}_+^{m+s} : x_j \geq \sum_{i=1}^n \lambda_i x_{ij}, j=1, \dots, m, y_r \leq \sum_{i=1}^n \lambda_i y_{ir}, r=1, \dots, s, \sum_{i=1}^n \lambda_i = 1, \lambda_i \in \{0,1\}, i=1, \dots, n \right\}.$$

Ejercicio 4. Escribe el modelo DEA bajo su formulación de **multiplicadores**, rendimientos **constantes** a escala y orientación **input** para la **observación B** del *toy dataset A*.

Luego, obtén el modelo DEA bajo su formulación **envolvente**, rendimientos **constantes** a escala y orientación **input**.

Para ello, considera la siguiente tabla de Tucker:

MAXIMIZACIÓN		MINIMIZACIÓN	
RESTRICCIONES	\leq	VARIABLES	\geq
	\geq		\leq
	$=$		$><$
VARIABLES	\geq	RESTRICCIONES	\geq
	\leq		\leq
	$><$		$=$

Ejercicio 5. Determina los siguientes modelos en formato matricial:

- Modelo DEA bajo su formulación de **multiplicadores**, rendimientos **constantes** a escala y orientación **input** para la **observación B** del *toy dataset A*:
- Modelo DEA bajo su formulación **envolvente**, rendimientos **constantes** a escala y orientación **input** para la **observación B** del *toy dataset A*:

Preguntas de contenido práctico

La base de datos `rice` contiene una muestra de 344 observaciones (la evolución de 43 productores de arroz en la región de Tarlac durante 8 años). Para el análisis, se han registrado las siguientes variables:

- `area`: área plantada en hectáreas (*input*).
- `labor`: días laborales de trabajo + días de trabajo contratado (*input*).
- `npk`: fertilizantes medidos en kilogramos de ingredientes activos (*input*).
- `prod`: toneladas de arroz producido (*output*).

Aunque esta base de datos refleja la evolución temporal de 43 productoras de arroz, a modo ilustrativo, las trataremos como 344 DMUs diferentes.

Por otro lado, también se utilizarán los *toy datasets* A, B, C y D.

Ejercicio 6. Crea la función `get_multipliers()` que devuelva los vectores de los multiplicadores (en formato `data.frame`) bajo una tecnología con rendimientos constantes a escala y una medida de eficiencia radial introducida por el usuario. Los argumentos de la función serán los siguientes:

- `tech_xmat`: matriz de *inputs* para las DMUs que constituyen la tecnología.
- `tech_ymat`: matriz de *outputs* para las DMUs que constituyen la tecnología.
- `eval_xmat`: matriz de *inputs* para las DMUs que se desean evaluar.
- `eval_ymat`: matriz de *outputs* para las DMUs que se desean evaluar.
- `measure`:
 - a) `"rad_out"`: distancia radial bajo orientación *output*.
 - b) `"rad_inp"`: distancia radial bajo orientación *input*.
- `rownames`: vector que indica el nombre de las DMUs. Puede ser `NULL`.

Ejercicio 7. Obtén los multiplicadores para la base de datos `rice` y la medida radial `input`. Luego, responde a las cuestiones.

- a) ¿Qué 5 unidades tienen mayor dependencia del `input area` para ser determinados como eficientes?

DMU	area	labor	npk	prod	v_1	v_2	v_3	μ_1	θ

b) ¿Qué 5 unidades reparten más equitativamente el peso entre los 3 inputs?

DMU	area	labor	npk	prod	v_1	v_2	v_3	μ_1	θ

c) Inserta una matriz de correlaciones entre los pesos de los 3 inputs y el score de eficiencia utilizando la librería `corrplot` y comenta los resultados

Ejercicio 8. Crea una función `eff_scores()` que devuelva un vector de *scores* (en formato `data.frame`) a partir de ciertas características de la tecnología y una medida de eficiencia introducida por el usuario. Los argumentos de la función serán los siguientes:

- `tech_xmat`: matriz de *inputs* para las DMUs que constituyen la tecnología.
- `tech_ymat`: matriz de *outputs* para las DMUs que constituyen la tecnología.
- `eval_xmat`: matriz de *inputs* para las DMUs que se desean evaluar.
- `eval_ymat`: matriz de *outputs* para las DMUs que se desean evaluar.
- `convexity`
 - `TRUE`: se impone convexidad en la tecnología.
 - `FALSE`: no se impone convexidad en la tecnología.
- `returns`
 - `"constant"`: se asumen rendimientos constantes a escala.
 - `"variable"`: se asumen rendimientos variables a escala.

- **measure:**
 - "rad_out": distancia radial bajo orientación *output*.
 - "rad_inp": distancia radial bajo orientación *input*.
 - "ddf": distancia direccional.

El argumento **direction** sólo se refiere a esta medida:

- "mean": vector proyección dado por el valor promedio de *inputs* y *outputs* de todas las DMUs.
 - "briec": vector proyección dado por el valor de *inputs* y *outputs* de la DMU evaluada.
- **rownames:** vector que indica el nombre de las DMUs. Puede ser NULL.

Ejercicio 9. Representa las siguientes fronteras de producción estimadas con la ayuda de **ggplot2**:

Nombre	Dataset	Convexidad	Libre disponibilidad	Rendimientos
CCR	<i>Toy dataset A</i>	✓	✓	Constantes
FDH	<i>Toy dataset B</i>	✗	✓	Variables

Frontera DEA (CCR) para el *toy dataset A*

Frontera FDH para el *toy dataset B*

Ejercicio 10. Realiza un análisis de la eficiencia para la base de datos **rice**:

- Obtén los *scores* de eficiencia mediante los modelos **DEA (CCR)**, **DEA (BCC)** y **FDH** bajo orientación radial *input* y *output*.
- Luego, obtén un gráfico comparativo de las densidades de los *scores* bajo orientación radial *input* obtenidos por los diferentes modelos e interprétalo.

- c) Interpreta los resultados de la DMU 4 (para el primer año) mediante el modelo DEA (BCC) con las 4 medidas que se indican:

DMU	area	labor	npk	prod	BCC.I	BCC.O	DDF.B	DDF.M
4	1.4	68	88	4.83	0.59	1.57	0.26	0.15
					Medida radial <i>input</i>			
					Medida radial <i>output</i>			
					DDF: briec			
					DDF: mean			

Ejercicio 11. Utiliza la base de datos D para responder las siguientes cuestiones:

- [1]. Crea una variable categórica binaria (t_{eff}) que determine si la DMU evaluada es verdaderamente eficiente con respecto a la frontera teórica de producción.
- [2]. Crea una nueva variable categórica binaria (r_{eff}) que determine si la DMU es eficiente en términos relativos a la muestra utilizando cualquier medida de eficiencia.
- [3]. ¿Cuántas DMU hay verdaderamente eficientes? En este sentido, ¿Qué quiere decir que DEA determina eficiencias relativas?
- [4]. Haz un gráfico que represente:
 - (I) La muestra observada de DMUs mediante un gráfico de dispersión.
 - (II) La frontera teórica (negro).
 - (III) La frontera aproximada obtenida por una modelo DEA (BCC) ("#DA33FF").
 - (IV) Un filtro de forma que permita identificar a las DMUs verdaderamente eficientes.
 - (V) Un filtro de color que permita identificar a las DMUs "relativamente" eficientes.
- [5]. Representa los vectores unitarios con el que se proyectará la DMU 23 bajo la medida de función de distancia direccional:
 - Color: "#33FF86" para la proyección bajo $direction = "mean"$.
 - Color: "#A6B209" para la proyección bajo $direction = "briec"$.
- [6]. Interpreta los resultados.

Anexo

Modelo **radial input** en formato de **ratio** con **tecnología CCR**:

$$\begin{aligned} \max_{\mathbf{v}, \boldsymbol{\mu}} \quad & \frac{\sum_{r=1}^s \mu_r \cdot y_{0r}}{\sum_{j=1}^m v_j \cdot x_{0j}} \\ \text{subject to} \quad & \frac{\sum_{r=1}^s \mu_r \cdot y_{ir}}{\sum_{j=1}^m v_j \cdot x_{ij}} \leq 1, \quad i = 1, \dots, n \\ & v_j \geq 0, \quad j = 1, \dots, m \\ & \mu_r \geq 0, \quad r = 1, \dots, s \end{aligned}$$

Modelo **radial output** en formato de **ratio** con **tecnología CCR**:

$$\begin{aligned} \min_{\mathbf{v}, \boldsymbol{\mu}} \quad & \frac{\sum_{j=1}^m v_j \cdot x_{0j}}{\sum_{r=1}^s \mu_r \cdot y_{0r}} \\ \text{subject to} \quad & \frac{\sum_{j=1}^m v_j \cdot x_{ij}}{\sum_{r=1}^s \mu_r \cdot y_{ir}} \geq 1, \quad i = 1, \dots, n \\ & v_j \geq 0, \quad j = 1, \dots, m \\ & \mu_r \geq 0, \quad r = 1, \dots, s \end{aligned}$$

Modelo **radial input** en formato de **multiplicadores** con tecnología DEA (CCR / BCC):

$$\max_{\mathbf{v}, \mathbf{\mu}, w} \sum_{r=1}^s \mu_r \cdot y_{0r} + w$$

subject to

$$\begin{aligned} \sum_{j=1}^m v_j \cdot x_{0j} &= 1 \\ \sum_{r=1}^s \mu_r \cdot y_{ir} + w &\leq \sum_{j=1}^m v_j \cdot x_{ij}, \quad i = 1, \dots, n \\ v_j &\geq 0, \quad j = 1, \dots, m \\ \mu_r &\geq 0, \quad r = 1, \dots, s \end{aligned}$$

- Tecnología **CCR**: $w = 0$.
- Tecnología **BCC**: w libre.

Modelo **radial output** en formato de **multiplicadores** con tecnología DEA (CCR / BCC):

$$\min_{\mathbf{v}, \mathbf{\mu}, w} \sum_{j=1}^m v_j \cdot x_{0j} + w$$

subject to

$$\begin{aligned} \sum_{r=1}^s \mu_r \cdot y_{0r} &= 1 \\ \sum_{j=1}^m v_j \cdot x_{ij} + w &\geq \sum_{r=1}^s \mu_r \cdot y_{ir}, \quad i = 1, \dots, n \\ v_j &\geq 0, \quad j = 1, \dots, m \\ \mu_r &\geq 0, \quad r = 1, \dots, s \end{aligned}$$

- Tecnología **CCR**: $w = 0$.
- Tecnología **BCC**: w libre.

Modelo **radial input** en formato **envolvente** con **tecnología CCR**:

$$\min_{\theta, \lambda} \theta$$

subject to

$$\sum_{i=1}^n \lambda_i \cdot x_{ij} \leq \theta x_{0j} \quad j = 1, \dots, m$$

$$\sum_{i=1}^n \lambda_i \cdot y_{ir} \geq y_{0r}, \quad r = 1, \dots, s$$

$$\lambda_i \geq 0, \quad i = 1, \dots, n$$

Tecnología **BCC**:

- Añadimos la restricción: $\sum_{i=1}^n \lambda_i = 1$

Tecnología **FDH**:

- Añadimos la restricción: $\sum_{i=1}^n \lambda_i = 1$
- Añadimos la restricción: $\lambda_i \in \{0,1\}, i = 1, \dots, n$.

Modelo **radial output** en formato **envolvente** con **tecnología CCR**:

$$\max_{\phi, \lambda} \phi$$

subject to

$$\sum_{i=1}^n \lambda_i \cdot x_{ij} \leq x_{0j} \quad j = 1, \dots, m$$

$$\sum_{i=1}^n \lambda_i \cdot y_{ir} \geq \phi y_{0r}, \quad r = 1, \dots, s$$

$$\lambda_i \geq 0, \quad i = 1, \dots, n$$

Tecnología **BCC**:

- Añadimos la restricción: $\sum_{i=1}^n \lambda_i = 1$

Tecnología **FDH**:

- Añadimos la restricción: $\sum_{i=1}^n \lambda_i = 1$
- Añadimos la restricción: $\lambda_i \in \{0,1\}, i = 1, \dots, n$.

Modelo **función distancia direccional** en formato **envolvente** con **tecnología CCR**:

$$\max_{\beta, \lambda} \beta$$

subject to

$$\sum_{i=1}^n \lambda_i \cdot x_{ij} \leq x_{0j} - \beta G_{x_j} \quad j=1, \dots, m$$

$$\sum_{i=1}^n \lambda_i \cdot y_{ir} \geq y_{0r} + \beta G_{y_r}, \quad r=1, \dots, s$$

$$\lambda_i \geq 0, \quad i=1, \dots, n$$

direction = “mean”: $G_{x_j} = \bar{x}_j$ y $G_{y_r} = \bar{y}_r$.

direction = “briec”: $G_{x_j} = x_{0j}$ y $G_{y_r} = y_{0r}$.

Tecnología **BCC**:

- Añadimos la restricción: $\sum_{i=1}^n \lambda_i = 1$

Tecnología **FDH**:

- Añadimos la restricción: $\sum_{i=1}^n \lambda_i = 1$
- Añadimos la restricción: $\lambda_i \in \{0,1\}, i=1, \dots, n$.