

# Propensity Score Matching

## Version Compacta

Dr. Sergio Bejar Lopez

CIDE

Enero 2026

# El Problema: Auto-Selección

## Ejemplo: Jóvenes Construyendo el Futuro

Programa de capacitación laboral - participación voluntaria

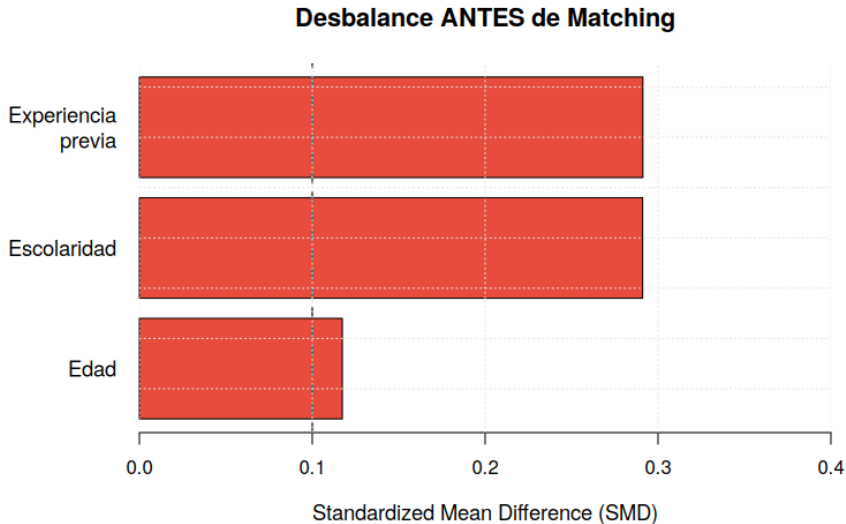
### Comparación ingenua:

Grupo	Empleo Formal
Participantes	44%
No participantes	38%
Diferencia	6 pp

### Problema

¿Efecto del programa o diferencias pre-existentes?

# Los Grupos NO Son Comparables



# La Solucion: Matching

## Idea central:

### Matching

Comparar cada participante con no-participantes **similares** en características observables

## Analogia:

- Encontrar "gemelos estadísticos"
- Uno participa, otro no
- Comparar estos pares

## Resultado:

- Grupos más comparables
- Estimación con menos sesgo
- Controla por observables

# Propensity Score: El Resumen

**Problema:** Muchas variables (edad, educacion, experiencia...)

**Solucion: Propensity Score**

## Definicion

$$e(X_i) = P(\text{Tratamiento} = 1|X_i)$$

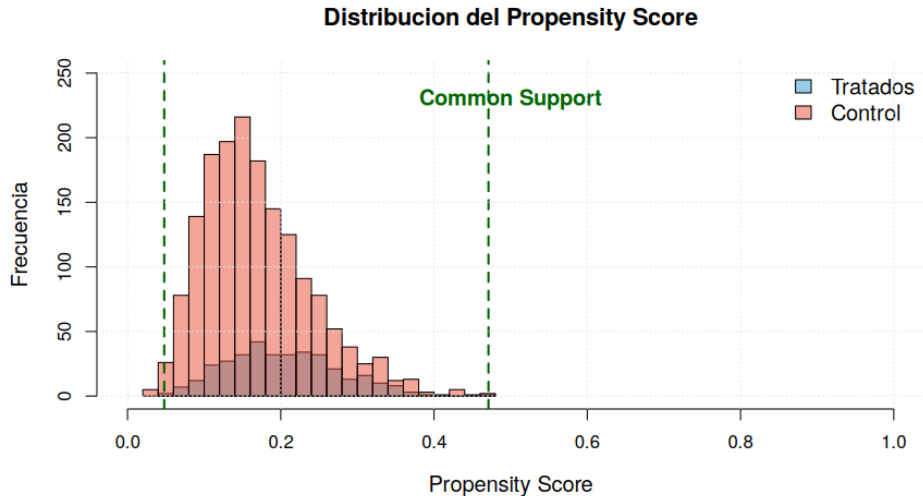
Probabilidad de participar dadas tus características

## Ventaja:

- Resume TODA la info en 1 numero (entre 0 y 1)
- Matching en PS = matching en todas las variables
- Facil de visualizar

**Calculo:** Regresion logistica

# Supuesto 1: Common Support



## Supuesto 2: CIA

### Conditional Independence Assumption (CIA)

$$\{Y_i(1), Y_i(0)\} \perp T_i | X_i$$

Condiciona en  $X$ , el tratamiento es como aleatorio

#### En palabras:

- Controlando por observables (edad, educación...)
- NO hay confusores no observados importantes
- Motivación, talento, redes no medidas pueden violar CIA

### CRITICO

CIA NO es testeable - requiere conocimiento del contexto

Si hay selección fuerte en no-observables: matching NO funciona

# Matching en R: Paso a Paso

```
library(MatchIt)
library(cobalt)

# Paso 1: Hacer matching
match <- matchit(
  tratamiento ~ edad + mujer + escolaridad + urbano +
    experiencia_previa + ingreso_familiar,
  data = datos,
  method = "nearest",
  distance = "glm"
)

# Paso 2: Verificar balance
bal.tab(match, thresholds = c(m = 0.1))

# Paso 3: Extraer datos matched
datos_matched <- match.data(match)

# Paso 4: Estimar ATT
att <- mean(datos_matched$empleo_formal[datos_matched$tratamiento==1]) -
  mean(datos_matched$empleo_formal[datos_matched$tratamiento==0])
```



# Balance: Antes vs Despues

## ANTES de matching (desbalanceado):

- Edad:  $SMD = 0.31$
- Escolaridad:  $SMD = 0.28$
- Experiencia:  $SMD = 0.19$

## DESPUES de matching (balanceado):

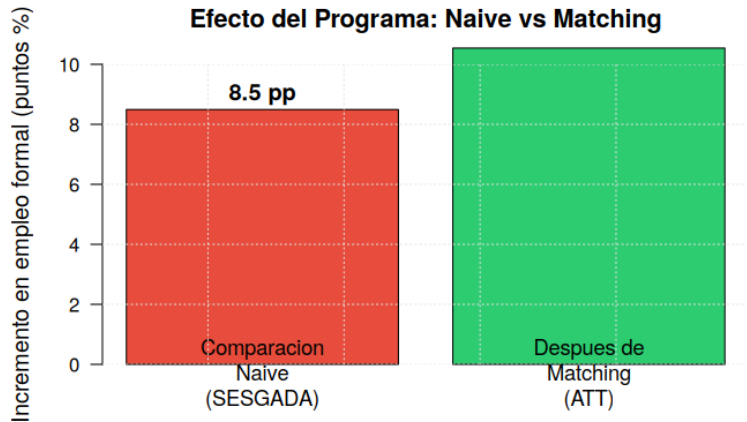
- Edad:  $SMD = 0.02$  ✓
- Escolaridad:  $SMD = 0.03$  ✓
- Experiencia:  $SMD = 0.01$  ✓

Exito!

$SMD < 0.1$  para todas las variables

Los grupos ahora son comparables

# Resultados: Naive vs Matching



# Limitaciones de Matching

## Ventajas:

- Intuitivo
- Transparente (podemos verificar balance)
- No requiere aleatorizacion

## Limitaciones IMPORTANTES:

- ❶ **CIA es fuerte** - si hay confusores no observados, falla
- ❷ **Solo controla observables** - lo no medido sesga
- ❸ **Requiere overlap** - sin common support no funciona
- ❹ **Pierde datos** - controles no matched se descartan

## Pregunta clave

¿Tengo medidas de TODAS las variables confusoras?

Si NO, considera DiD, RDD, o IV

# Matching vs DiD

	Matching	DiD
Datos	Cross-section	Panel
Supuesto	CIA	Tendencias paralelas
Controla	Observables	Fijos + tendencias
Ventaja	No necesita panel	Controla no-observables

## Combinables

DiD + Matching = robusto a ambos problemas

# Resumen: 4 Puntos Clave

- 1 **Matching** = comparar similares en observables
- 2 **Propensity Score** = resume info en 1 numero
- 3 **Supuestos:** CIA + Common Support
- 4 **Verificar balance:**  $SMD < 0.1$

## Mensaje final

Matching es poderoso PERO solo si CIA es plausible

# Codigo Completo

```
1 library(MatchIt)
2 library(cobalt)
3
4 # Matching
5 m <- matchit(tratamiento ~ edad + mujer + escolaridad +
6             urbano + experiencia_previa,
7             data = datos,
8             method = "nearest")
9
10 # Balance
11 bal.tab(m, thresholds = c(m = 0.1))
12
13 # Datos matched
14 datos_m <- match.data(m)
15
16 # ATT
17 att <- mean(datos_m$outcome[datos_m$tratamiento==1]) -
18        mean(datos_m$outcome[datos_m$tratamiento==0])
19
20 # Regresion (mas robusto)
21 modelo <- lm(outcome ~ tratamiento,
22             data = datos_m,
23             weights = weights)
24 summary(modelo)
```

# Matching

Comparar Similares

Dr. Sergio Bejar Lopez  
CIDE