# IA Responsable — Proyecto 01

Este sitio resume el proyecto de **IA Responsable** aplicando los principios **FATE** (Fairness, Accountability, Transparency, Ethics) sobre el dataset Adult Census Income (UCI). Se entrena un modelo base y una versión mitigada usando class\_weight="balanced", y se evalúa el desempeño **segmentado por grupos sensibles** (sexo, edad y raza).

- Instrucciones de reproducción: ver README.md
- Código completo: project-eda.ipynb · baseModel.ipynb

- Dataset: 32,561 filas · 15 columnas · atributos sensibles usados para segmentar: sex , race , Mejor baseline (sin balanceo) — Logistic Regression: Accuracy 0.8096, F1 0.5110, ROC AUC
- 0.8221. Mitigación (class\_weight="balanced") — Logistic Regression: Accuracy 0.7375, F1 0.5797,
- Efecto: la mitigación aumenta el recall/TPR de la clase positiva (FN ↓, TP ↑), mejora F1 y reduce algunas brechas de Equal Opportunity (TPR por grupo), a costa de mayor FPR y Positive Rate
- en varios subgrupos.
- 2) Evaluar el impacto por subgrupos sensibles,

Objetivo del proyecto: 1) Identificar posibles sesgos en datos y modelo,

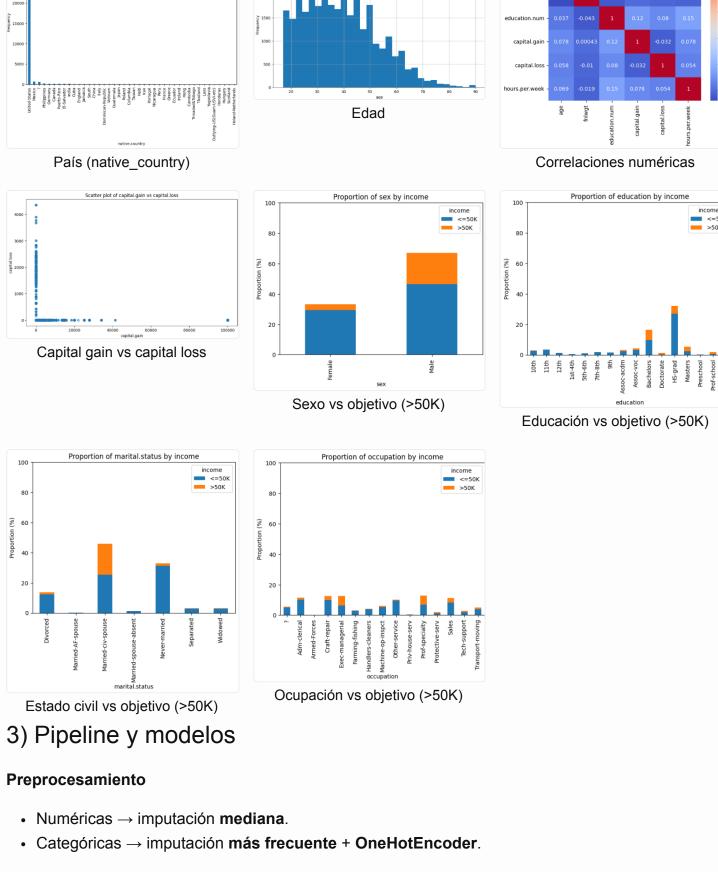
### 3) Aplicar una mitigación y comparar antes vs. después, 4) **Reflexionar** sobre riesgos, transparencia y responsabilidad.

- 2) Datos (Adult Census Income UCI)

#### • Tamaño: 32,561 registros, 15 columnas (6 numéricas, 9 categóricas). • Valores "?": workclass (1,836), occupation (1,843), native.country (583). Atributos sensibles: sex , race y age (agrupada en bins).

- El EDA (ver notebook) reporta distribuciones por categoría, estadísticas numéricas y manejo de "?" como NaN y posterior imputación/One-Hot.
- desbalance moderado.
- Reemplazo de "?" por NaN en categóricas antes de estadísticas/plots.
- Imputación: mediana (numéricas) y moda (categóricas). vistas.

Other Amer-Indian-Eskimo



## Modelo Logistic (baseline)

Tree (baseline)

Logistic (balanced)

4) Métricas globales (Test)

Baseline: Logistic Regression, Decision Tree.

Modelos evaluados

**Partición** 

0.6593 0.5473 Tree (balanced)

• Mitigación: mismas arquitecturas con class\_weight="balanced" .

• Train 75% / Test 25%, stratify=y, random state=42.

Lectura La versión **balanced** sacrifica *accuracy* (más FP) para **recuperar muchos positivos** (menos FN) → sube F1. El AUC (ranking) se mantiene ≈ igual, consistente con cambiar la penalización/umbral más que la capacidad discriminativa. 5) Equidad: evaluación segmentada por subgrupos Se reportan métricas por grupo: accuracy, F1, TPR (recall), FPR, Positive Rate. También se muestran "brechas" (max-min entre grupos) para cada métrica: menor es mejor.

F1

0.5549

0.1333

0.6130

**TPR** 

0.4730

0.0748

0.8175

gap\_tpr

0.3982

0.4366

**FPR** 

0.1018

0.0059

**FPR** 

0.3734

gap\_fpr

0.0959

0.2751

Accuracy

0.8096

0.7978

0.7375

F1

0.5110

0.4607

0.5797

**ROC AUC** 

0.8221

0.8029

0.8231

0.8061

**Positive Rate** 

0.2151

0.0134

**Positive Rate** 

0.5090

0.1293

gap\_positive\_rate

0.2017

0.3796

**Positive Rate** 

0.4298

0.1064

0.6430

**Positive Rate** 

0.4127

0.1490

0.3891

0.1772

0.1111

gap\_positive\_rate

0.1498

0.3016

#### Después (Logistic — balanced) accuracy F1 **TPR** sex n

0.6849

accuracy

0.7684

0.8934

5.1. Sexo ( sex ) — antes vs. después

n

5458

2683

5458

gap\_accuracy

0.1250

0.1597

Antes (Logistic — baseline)

sex

Male

Male

Modelo

Baseline

**Balanced** 

Después (Logistic — balanced)

3927

2407

1647

personas mayores— y la tasa de positivos predicha.

7017

725

257

79

63

gap\_accuracy

0.1204

0.1687

5.3. Raza (race) — antes vs. después

Antes (Logistic — baseline)

race

White

Black

Other

Modelo

Baseline

**Balanced** 

Asian-Pac-Islander

Amer-Indian-Eskimo

Brechas (max-min)

age\_bin

30-49

<30

50-69

Female

<b>Resumen</b> : la mitigación <b>eleva TPR</b> y <b>mejora F1</b> en ambos sexos (mejor <i>Equal Opportunity</i> ), pero <b>aumenta FPR y Positive Rate</b> , ampliando brechas en <i>Demographic Parity</i> .									
5.2. Edad (age, bins) — antes vs. después									
Antes (Logistic — baseline)									
age	bin	n	accuracy	F1	TPR	FPR	Positive Rate		

#### 0.4062 0.4025 0.8649 70+ 160 0.7317 0.7625 Brechas (max-min)

F1

0.5979

0.2674

0.6387

accuracy

0.7041

0.8816

0.6387

**TPR** 

0.7188

0.3910

0.8946

**FPR** 

0.3024

0.0897

0.5033

White	7017	0.8002	0.5212	0.4253	0.0710	0.1616
Black	725	0.8855	0.2783	0.1818	0.0173	0.0372
Asian-Pac-Islander	257	0.8093	0.5586	0.4697	0.0733	0.1751
Amer-Indian-Eskimo	79	0.8608	0.0000	0.0000	0.0286	0.0253
Other	63	0.9206	0.0000	0.0000	0.0333	0.0317

Resumen: mejora el IPR (oportunidad) para varios grupos, pero incrementan FPR y Positive Rate; en								
grupos con muy bajo soporte (por ejemplo, Other y Amer-Indian-Eskimo) persisten								
nestabilidades.								
6) Mitigación aplicada y alternativas								
Aplicada: class_weight="balanced" en Logistic y Tree.								
<ul> <li>Ventajas: ↑ TPR/recall de la clase positiva, ↑ F1, ↓ brechas EO en varios casos.</li> </ul>								
<ul> <li>Costos: ↑ FPR y ↑ Positive Rate → posibles efectos en Demographic Parity.</li> </ul>								
Alternativas recomendadas								
Ajuste de umbral por grupo (optimización precision–recall con restricciones de EO/EqOdds).								
<ul> <li>Rebalanceo de datos (over/under-sampling) o reweighing previo.</li> </ul>								
Post-procesamiento (p. ej., Equalized Odds).								
<ul> <li>Agrupar/reescalar categorías con n muy bajo para reducir varianza.</li> </ul>								

# • Métricas de equidad: definir qué se reporta (TPR/EO, FPR, Positive Rate/DP) y por qué. • Publicación: notebooks y scripts reproducibles; requirements.txt y pasos en README. **Ética (Ethics)**

producción sin justificación.

históricamente subrepresentados.

Riesgos al desplegar en entornos reales

de ranking no se deteriora).

Limitaciones

F1/TPR).

fluctuar.

Responsabilidad (Accountability)

7) Reproducibilidad

python -m pip install -r requirements.txt

Requisitos

• Umbral: pequeñas variaciones alteran TPR/FPR; conviene análisis de sensibilidad. 9) Conclusiones y recomendaciones Qué mejoró (equidad y desempeño)

• Label/selection bias: etiquetas históricas pueden arrastrar sesgos sistémicos.

• Distribution shift: el rendimiento y la equidad pueden cambiar si cambia la población.

Selección de criterios de equidad: justificar EO/DP según el caso de uso.

baseline. Qué no mejoró / trade-offs observados

# desproporcionadas. etiquetas/datos). Shift de distribución: cambios en la población degradan desempeño y equidad si no hay

• Gobernanza y cumplimiento: se requiere trazabilidad, explicaciones y mecanismos de apelación para usuarios afectados.

- Post-procesamiento de fairness (p. ej., equalized odds/EO) para controlar TPR/FPR por grupo sin reentrenar todo el pipeline.
- de **drift** y auditorías periódicas. • Proceso humano-en-el-bucle: revisión de casos límite, documentación (model card), registro de decisiones y canal de apelación.

Repositorio: Enlace

- TL;DR Resultados clave
  - **ROC AUC 0.8231.**
- 1) Contexto y objetivo Los modelos de ML pueden heredar o amplificar sesgos presentes en los datos o en el diseño/uso del sistema.
- Hallazgos clave del EDA
- Distribución del target: ~24% de la clase positiva (>50K) y ~76% de la negativa (<=50K) →</li> Workclass: predominio del sector privado (Private) frente a gobierno/cuenta propia → posible sesgo por sobrerrepresentación.
- Educación: concentración entre secundaria—universidad; menor presencia en niveles extremos.
- Raza: ~85% White; otras razas con bajo soporte → riesgo de alta varianza e inequidad si no se • País (native country): mayoría Estados Unidos; múltiples categorías raras con muy pocos • Edad: mayor densidad entre 20-50 años (población laboral activa); útil estratificar en bins para análisis. Tratamiento de datos (según EDA/Pipeline)
- Codificación: OneHotEncoder(handle\_unknown="ignore") para robustez ante categorías no Variables candidatas usadas • Numéricas: age , education\_num , hours\_per\_week , capital\_gain , capital\_loss . • Categóricas: workclass, education, marital status, occupation, sex, race, native country. Gráficas (EDA) Proportion of workclass

# Educación Raza Trabajo (workclass)

# Matrices de confusión (test) Logistic — baseline: TN=5781, FP=400, FN=1150, TP=810 Logistic — balanced: TN=4530, FP=1651, FN=486, TP=1474

#### Female 2683 0.8446 0.3495 0.3810 0.0984 **Brechas (max-min)**

gap\_f1

0.4216

0.2635

30–49	3927	0.7601	0.4875	0.3727	0.0690	0.1620
<30	2407	0.9393	0.1412	0.0902	0.0110	0.0154
50–69	1647	0.7511	0.6103	0.5459	0.1350	0.2817
70+	160	0.6750	0.5273	0.7838	0.3577	0.4562

#### Modelo gap\_positive\_rate gap\_accuracy gap\_f1 gap\_tpr gap\_fpr 0.2643 0.6936 0.3467 0.4409 Baseline 0.4691 Balanced 0.4753 0.5036 0.6420 0.6561 0.3714

accuracy

0.7241

0.8566

0.7043

0.8354

0.8730

gap\_f1

0.5586

0.3872

Resumen: la mitigación sube TPR en todos los bins (recupera positivos), pero aumenta FPR —más en

F1

**TPR** 

0.7676

0.5227

0.6818

0.5556

0.3333

gap\_fpr

0.0560

0.1935

**FPR** 

0.2908

0.0973

0.2880

0.1286

0.1000

race	n	accuracy	F1	TPR	FPR	Positive Rate		
Después (Logistic — balanced)								
Other	63	0.9206	0.0000	0.0000	0.0333	0.0317		
Amer-Indian-Eskimo	79	0.8608	0.0000	0.0000	0.0286	0.0253		
Asian-Pac-Islander	257	0.8093	0.5586	0.4697	0.0733	0.1751		

0.5872

0.4694

0.5422

0.4348

0.2000

gap\_tpr

0.4697

0.4342

 Revisión por pares: documentar decisiones (por qué este umbral/mitigación) y sus efectos. Transparencia (Transparency) • Pipeline documentado: imputación, codificación, modelo, umbral de decisión y atributos sensibles.

• Riesgos y uso responsable: evitar despliegues sin monitoreo; impacto en grupos vulnerables.

Dataset: desbalance de clases; subgrupos con n muy bajo → alta varianza (inestabilidades en

• ↑ TPR (Equal Opportunity) en la mayoría de subgrupos con class weight="balanced" → se reducen falsos negativos y aumenta la probabilidad de acierto para la clase positiva en grupos

↑ F1 de la clase positiva gracias al aumento del recall, manteniendo ROC AUC similar (la capacidad

• Gobernanza de datos: respetar licencias y privacidad; no introducir atributos sensibles en

Versionado y trazabilidad: commits atómicos con mensajes claros; notebooks con random state

Registro de experimentos: anotar fecha, datos, hiperparámetros, métricas y figuras exportadas.

8) Responsabilidad, transparencia, ética y limitaciones

 ↑ FPR y ↑ Positive Rate en varios subgrupos → puede ampliar brechas de Demographic Parity (más positivos predichos en ciertos grupos).

• En varias segmentaciones (sexo/raza), la brecha de TPR (gap EO) se estrecha respecto al

• Impacto desigual: mayores FPR en ciertos grupos pueden traducirse en costos o fricciones Feedback loops: decisiones del modelo pueden reforzar sesgos (p. ej., si influyen futuras

 
 ↓ Accuracy global por el aumento de falsos positivos (costo esperado al priorizar recall/equidad).
 • Persisten inestabilidades en grupos con **muy bajo soporte** (n pequeño), donde F1/TPR pueden

Recomendaciones prácticas Ajuste de umbral (global y/o por grupo) guiado por curvas PR/ROC, priorizando objetivos de EO o **EqOdds** según el caso de uso. • Mitigaciones de datos: reweighing, sobre/sub-muestreo estratificado y ampliación de datos en

• Calibración por grupo (reliability curves/Platt/Isotónica) para alinear probabilidades y reducir sesgos de decisión. Monitoreo continuo: dashboard de métricas por subgrupo (TPR/FPR/PositiveRate, gaps), alertas

Equidad, responsabilidad, transparencia y ética aplicadas a un modelo predictivo con Adult Census.

IA Responsable — Proyecto 01 IA Responsable — Proyecto 01

subgrupos con **n** bajo.