

Proyecto Big Data - FULL DATA

Selección del Dataset asignado

Usamos el dataset Gym Members Exercise Tracking, que tiene 973 registros y 15 variables. Este dataset contiene datos de personas que entrenan en el gimnasio, con las siguientes variables: Age, Gender, Weight (kg), Height (m), Max_BPM, Avg_BPM, Resting_BPM, Session_Duration (hours), Calories_Burned, Workout_Type, Fat_Percentage, Water_Intake (liters), Workout_Frequency (days/week), Experience_Level y BMI.

Definición del problema

Queremos predecir el Experience_Level de un miembro del gimnasio (principiante, intermedio o avanzado) a partir de sus hábitos y características físicas.

Este análisis es importante porque en muchos gimnasios, al momento de la inscripción, no se conoce con precisión el nivel de experiencia real del nuevo miembro. Esto puede llevar a asignarle rutinas demasiado exigentes o, por el contrario, demasiado simples, generando frustración, desmotivación y un alto riesgo de abandono.

Contar con un modelo que prediga el nivel de experiencia permite a la empresa personalizar los planes de entrenamiento desde el primer día, ofrecer una mejor experiencia inicial y aumentar la retención de clientes. Además, esta información puede utilizarse para diseñar programas de fidelización más efectivos y optimizar la asignación de entrenadores y recursos dentro del gimnasio.

Análisis Exploratorio de Datos:

Las hipótesis que planteamos son las siguientes:

- **Hipótesis 1:** La variable Workout_Frequency (days/week) y la Session_Duration (hours) influyen directamente en el nivel de experiencia.
- **Hipótesis 2:** El Fat_Percentage y el BMI permiten diferenciar a principiantes de avanzados, ya que los principiantes suelen tener valores más altos y los avanzados más controlados.
- **Hipótesis 3:** El Workout_Type (Cardio, HIIT, Strength, Yoga) es un buen predictor del nivel de experiencia, porque los más avanzados tienden a llevar a cabo y combinar rutinas más intensas.

Antes de modelar, realizamos un análisis exploratorio para entender el comportamiento general del dataset y detectar patrones relevantes.

Primero observamos que la mayoría de los registros corresponden a principiantes e intermedios, mientras que los usuarios avanzados son minoría. Esto implica que el modelo podría tener más ejemplos de los primeros niveles, algo a tener en cuenta al evaluar su rendimiento.

También identificamos que las variables Workout_Frequency y Session_Duration son las que más diferencian a los grupos. Los avanzados entrenan más días por semana y durante más tiempo, mientras que los principiantes tienen rutinas más cortas y menos frecuentes. Este patrón confirma que la constancia y el tiempo de dedicación son factores claves en la progresión del usuario.

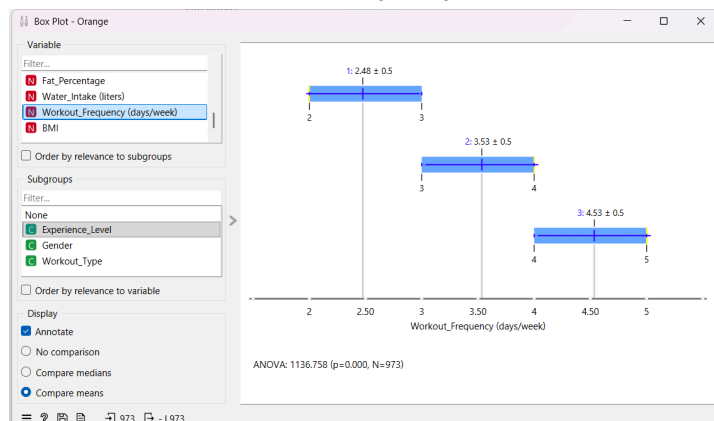
En cuanto a la composición corporal, los principiantes presentan un mayor porcentaje de grasa corporal, mientras que los avanzados muestran valores más bajos y estables. Sin embargo, el BMI no evidencia grandes diferencias, lo que sugiere que no es un indicador tan confiable del nivel de experiencia por sí solo.

Desde una mirada práctica, estos hallazgos ofrecen insights relevantes para una empresa de fitness:

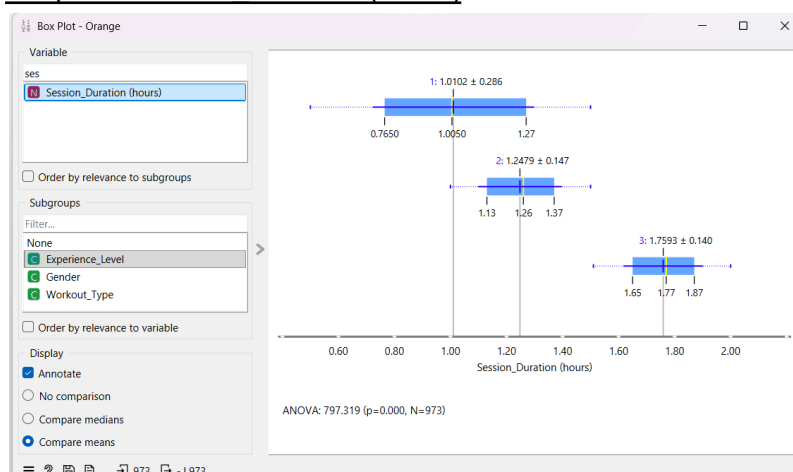
- Un nuevo miembro con baja frecuencia y sesiones cortas podría necesitar mayor acompañamiento para evitar el abandono.
- Observar el progreso en la duración o frecuencia de las rutinas puede servir como indicador temprano de fidelización.
- Las métricas corporales deben complementarse con hábitos de entrenamiento para definir mejor el nivel y las recomendaciones personalizadas.;

En resumen, el EDA no solo ayudó a validar nuestras hipótesis, sino también a entender qué comportamientos predicen mejor la constancia y la evolución del usuario, lo que es clave para mejorar la retención en gimnasios o apps de entrenamiento.

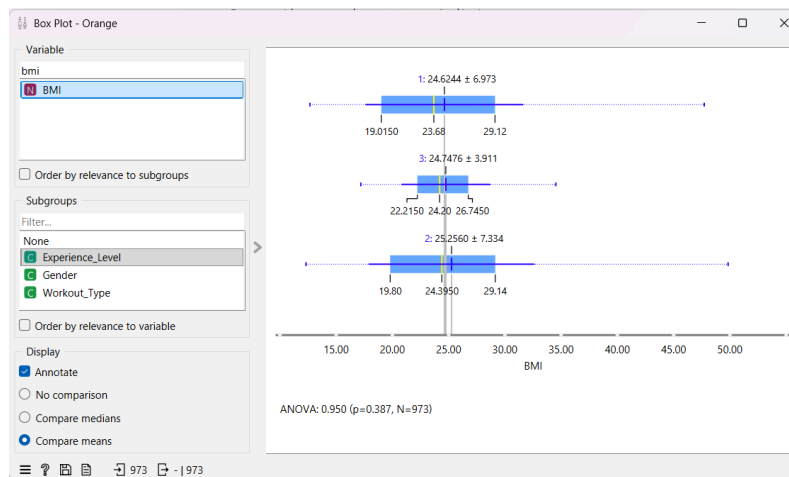
Boxplot - Workout Frequency (days/week)



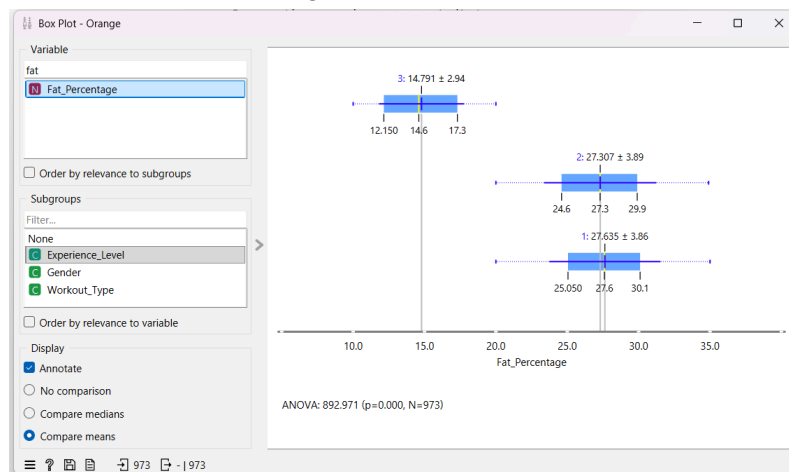
Boxplot - Session Duration (hours)



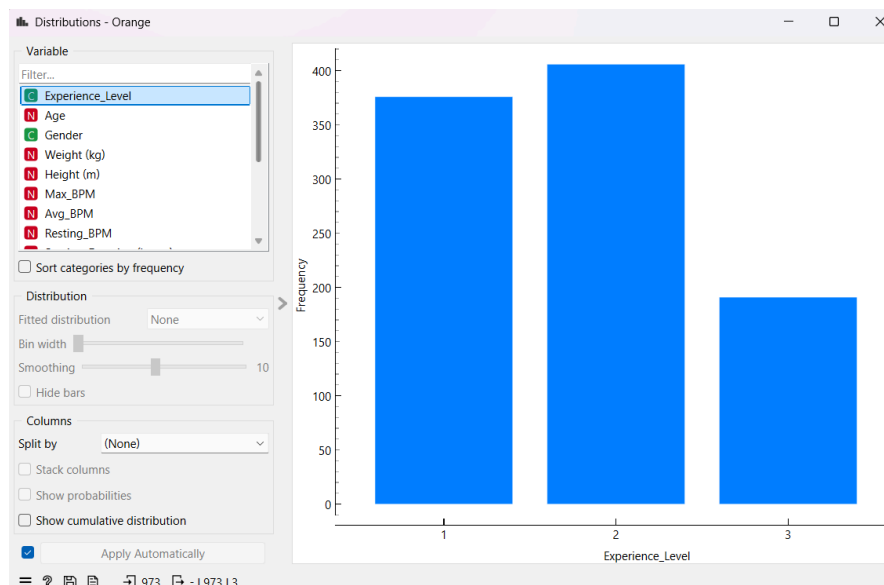
Boxplot - BMI



Boxplot - Fat_Percentage



Distribución de Experience_Level



Preprocesamiento y Selección de Variables

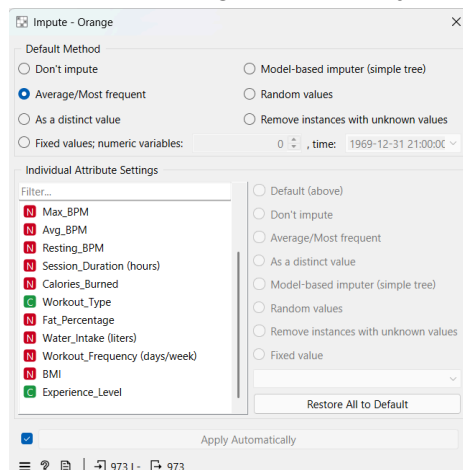
Antes de aplicar los modelos de clasificación, realizamos un preprocesamiento completo del dataset para garantizar la calidad de los datos y hacer más sencilla su interpretación para los algoritmos. En este proceso usamos distintos widgets de visualización como Data Table, Feature Statistics, Distributions y Box Plots, para ir verificando los cambios que fuimos haciendo en cada etapa y poder controlar la consistencia de los datos.

Select Columns

En el widget Select Columns establecimos la variable Experience_Level como Target. El resto de las variables se mantuvieron como Features, y no nos pareció necesario descartar ninguna, ya que todas aportaban información que podría resultar relevante para el modelo.

Impute

Aunque no registramos ningún valor nulo en el dataset, utilizamos por precaución el widget Impute con la configuración por defecto, para asegurar que no haya datos vacíos y mantener la integridad del conjunto.



Continuize

En este paso utilizamos dos widgets Continuize, cada uno con un propósito distinto dentro del flujo.

El primero se aplicó antes del Data Table (2), para preparar los datos del modelado. En este caso, realizamos un One-hot encoding sobre la variable Workout_Type, transformándola en columnas numéricas (una por cada tipo de entrenamiento) para que los modelos pudieran procesarla correctamente. La variable Gender se mantuvo categórica binaria, mientras que las variables numéricas conservaron la opción Keep as it is, sin normalización ni estandarización.

El segundo Continuize se utilizó solo para el análisis de correlaciones. Allí aplicamos un One-hot encoding a la variable Experience_Level, con el fin de convertirla en variables binarias (una por cada nivel de experiencia) y poder calcular las correlaciones entre esta y el resto de las variables.

Continue - Orange

Categorical Variables

- ★ Preset: first as base
- Filter...
- Gender
- Workout_Type
- Targets
- Experience_Level: one-hot

Numeric Variables

- ★ Preset: no change
- Filter...
- Age
- Weight (kg)
- Height (m)
- Max_BPM
- Avg_BPM
- Resting_BPM

Reset All

Apply Automatically

973

Continue (1) - Orange

Categorical Variables

- ★ Preset: first as base
- Filter...
- Gender
- Workout_Type: one-hot
- Targets
- Experience_Level

Numeric Variables

- ★ Preset: no change
- Filter...
- Age
- Weight (kg)
- Height (m)
- Max_BPM
- Avg_BPM
- Resting_BPM

Reset All

Apply Automatically

973

Data Table (2) - Orange

Info

973 instances (no missing data)

17 features

Target with 3 values

No meta attributes.

Variables

- ☒ Show variable labels (if present)
- ☐ Visualize numeric values
- ☒ Color by instance classes

Selection

- ☒ Select full rows

Restore Original Order

☒ Send Automatically

	Experience_Level	Age	Gender=Male	Weight (kg)	Height (m)
1	3	56	1	88.3	1.71
2	2	46	0	74.9	1.53
3	2	32	0	68.1	1.66
4	1	25	1	53.2	1.70
5	1	38	1	46.1	1.79
6	3	56	0	58.0	1.68
7	2	36	1	70.3	1.72
8	2	40	0	69.7	1.51
9	2	28	1	121.7	1.94
10	1	28	1	101.8	1.84
11	1	41	1	120.8	1.67
12	2	53	1	51.7	1.70
13	2	57	1	112.5	1.61
14	1	41	1	94.5	2.00
15	2	20	1	117.7	1.81

Reset All

Apply Automatically

973

Data Sampler

Usamos el widget de Data Sampler para dividir el dataset en dos subconjuntos: un 70% de los datos para entrenamiento y el otro 30% para prueba.

Data Sampler - Orange

Sampling Type

- ☒ Fixed proportion of data: 70 %
- ☐ Fixed sample size
- Instances: 2
- ☐ Sample with replacement
- ☐ Cross validation
- Number of subsets: 10
- Unused subset: 1
- ☐ Bootstrap

Options

- ☒ Replicable (deterministic) sampling
- ☐ Stratify sample (when possible)

Sample Data

973 | 682 | 291

Ya en esta instancia, el flujo de trabajo ya se encontraba preparado para avanzar con la etapa de modelado.

Modelado y Evaluación del Modelo

Para este caso, el error más costoso sería clasificar como avanzado a un usuario que en realidad es principiante, ya que recibiría rutinas demasiado exigentes y podría abandonar el gimnasio. Por eso, la métrica más importante para la empresa es la Precision, que mide la capacidad del modelo para evitar falsos positivos.

El modelo Random Forest fue el que mejor desempeño tuvo en este sentido, con una Precision de 0.897, lo que indica que logra identificar correctamente la mayoría de los usuarios avanzados sin confundir a principiantes. Además, mantiene un buen equilibrio con un Recall de 0.890 y un AUC de 0.953, mostrando una alta capacidad para distinguir entre los diferentes niveles de experiencia.

Por eso, seleccionamos Random Forest como el modelo más confiable para predecir el nivel de experiencia de los usuarios y apoyar decisiones que reduzcan el riesgo de abandono en los nuevos miembros del gimnasio.

Model	AUC	CA	F1	Prec	Recall
Logistic Regression	0.959	0.863	0.862	0.863	0.863
Naive Bayes	0.949	0.845	0.844	0.845	0.845
Random Forest	0.953	0.890	0.888	0.897	0.890

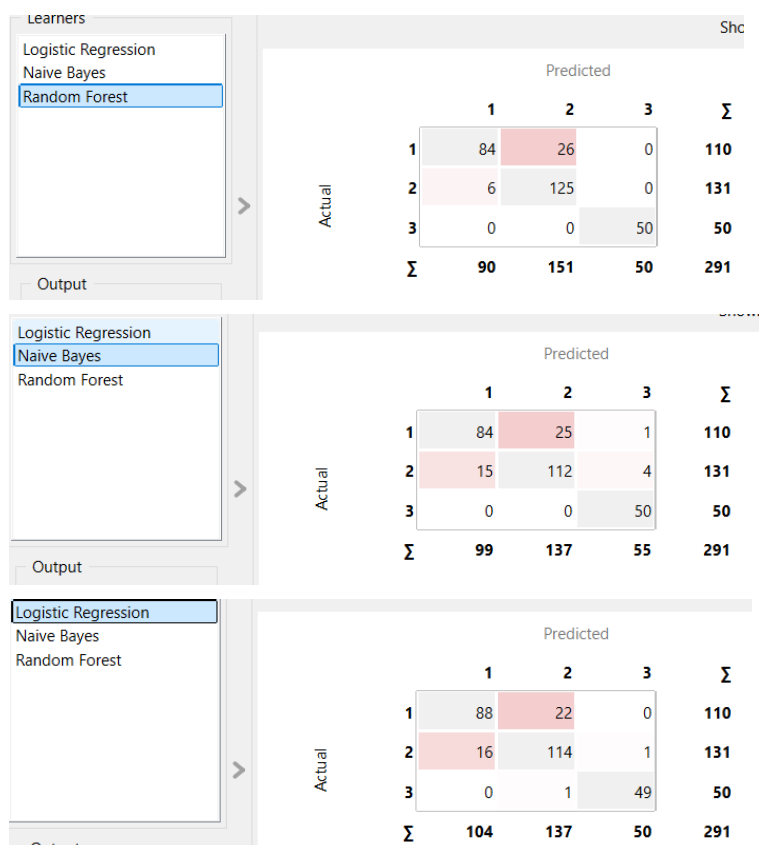
Matriz de Confusión

Al analizar las matrices de confusión, se observa que el modelo Random Forest logró la clasificación más equilibrada entre los distintos niveles de experiencia. La mayoría de los casos fueron correctamente identificados, y los errores se concentran principalmente entre clases contiguas (por ejemplo, confundir a un intermedio con un principiante), lo cual es esperable y menos grave.

Desde el punto de vista del negocio, el error más importante a evitar es clasificar como avanzado a alguien que no lo es, ya que podría recibir rutinas demasiado exigentes y abandonar el gimnasio. En ese sentido, Random Forest presenta el menor número de falsos positivos, lo que refuerza su alta Precision y su confiabilidad para predecir correctamente a los usuarios más experimentados sin sobreestimar su nivel.

En comparación, los modelos Naive Bayes y Logistic Regression muestran una mayor cantidad de errores cruzados entre niveles, lo que implica menor control sobre este tipo de equivocaciones.

Por eso, además de tener buenas métricas generales, Random Forest destaca por su distribución de errores más coherente y por minimizar los errores más costosos para la empresa.



Interpretación de las Predicciones

A partir de las correlaciones obtenidas con el modelo de Random Forest, se comprobó que las variables más influyentes para predecir el nivel de experiencia están relacionadas con la frecuencia, duración e intensidad del entrenamiento, tal como se anticipaba en las hipótesis iniciales.

Hipótesis 1:

La frecuencia de entrenamiento (Workout_Frequency) y la duración de las sesiones (Session_Duration) fueron las variables más relevantes. En los niveles iniciales se observan correlaciones negativas fuertes (alrededor de -0.8 y -0.7), lo que indica que quienes entrenan poco suelen ser principiantes. En cambio, a medida que aumenta la experiencia, estas correlaciones se vuelven positivas (en torno a $+0.6$ y $+0.7$), reflejando que los usuarios más avanzados entrenan con mayor constancia y por más tiempo.

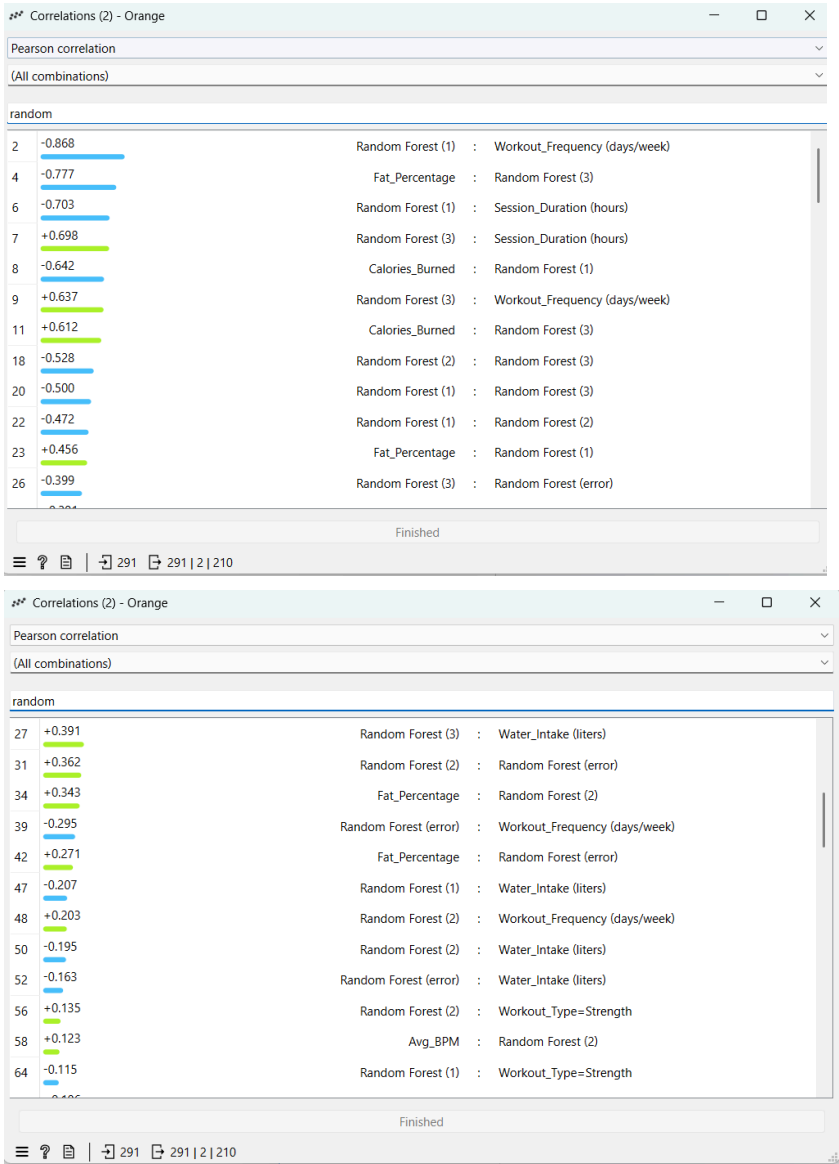
Hipótesis 2:

El porcentaje de grasa corporal (Fat_Percentage) también tuvo un papel destacado. Se registraron correlaciones negativas elevadas (cercanas a -0.75), lo que significa que a mayor nivel de experiencia, menor porcentaje de grasa corporal. Esto respalda la idea de que el progreso técnico suele ir acompañado de una mejora en la composición corporal.

Hipótesis 3:

El tipo de entrenamiento (Workout_Type) no mostró una relación marcada con el nivel de experiencia. Las correlaciones fueron débiles (por ejemplo, alrededor de -0.16 en el caso de Strength), lo que sugiere que tanto principiantes como avanzados suelen realizar tipos de ejercicio similares, y que el progreso depende más de la frecuencia, la duración y la regularidad que del tipo de rutina.

En conjunto, el modelo confirma que la frecuencia de entrenamiento, la duración de las sesiones y el control corporal son los factores que mejor explican el avance en el nivel de experiencia, coherentes con las tendencias detectadas en el análisis exploratorio previo.



Predictions - Orange																
Show probabilities for classes in data						Show classification errors					Restore Original Order					
	Logistic Regression		Naive Bayes		Random Forest					Experience_Level	Age	Gender=Male	Weight (kg)	Height (m)		
1	0.55 : 0.45 : 0.00 → 1	0.552	0.95 : 0.05 : 0.00 → 1	0.946	0.44 : 0.56 : 0.00 → 2	0.440				25	0		41.1	1.67	18	
2	0.37 : 0.63 : 0.00 → 2	0.635	0.40 : 0.60 : 0.00 → 2	0.604	0.48 : 0.52 : 0.00 → 2	0.518				31	0		42.7	1.76	18	
3	1.00 : 0.00 : 0.00 → 1	0.003	1.00 : 0.00 : 0.00 → 1	0.001	0.95 : 0.05 : 0.00 → 1	0.051				20	1		83.0	1.80	16	
4	0.25 : 0.75 : 0.00 → 2	0.250	0.40 : 0.60 : 0.00 → 2	0.397	0.45 : 0.54 : 0.01 → 2	0.458				50	1		96.7	1.72	18	
5	0.00 : 1.00 : 0.00 → 2	0.004	0.00 : 1.00 : 0.00 → 2	0.001	0.12 : 0.88 : 0.01 → 2	0.124				54	0		75.6	1.61	18	
6	0.00 : 0.00 : 1.00 → 3	0.000	0.00 : 0.00 : 1.00 → 3	0.001	0.00 : 0.03 : 0.97 → 3	0.030				21	1		88.4	1.60	19	
7	0.89 : 0.11 : 0.00 → 1	0.111	0.93 : 0.07 : 0.00 → 1	0.068	0.93 : 0.07 : 0.00 → 1	0.065				23	0		41.9	1.58	16	
8	0.01 : 0.99 : 0.01 → 2	0.010	0.00 : 1.00 : 0.00 → 2	0.004	0.07 : 0.93 : 0.00 → 2	0.071				55	1		63.5	1.86	19	
9	0.38 : 0.62 : 0.00 → 2	0.380	0.33 : 0.67 : 0.00 → 2	0.329	0.38 : 0.61 : 0.00 → 2	0.386				40	0		52.6	1.66	19	
10	0.57 : 0.43 : 0.00 → 1	0.572	0.53 : 0.47 : 0.00 → 1	0.526	0.48 : 0.52 : 0.00 → 2	0.476				28	0		62.4	1.60	17	
11	0.35 : 0.65 : 0.00 → 2	0.349	0.35 : 0.65 : 0.00 → 2	0.349	0.38 : 0.62 : 0.01 → 2	0.385				27	1		104.3	1.68	16	
12	0.00 : 0.00 : 1.00 → 3	0.000	0.00 : 0.00 : 1.00 → 3	0.000	0.00 : 0.00 : 1.00 → 3	0.000				42	1		85.2	1.81	18	
13	0.59 : 0.41 : 0.00 → 1	0.588	0.21 : 0.79 : 0.00 → 2	0.206	0.44 : 0.56 : 0.00 → 2	0.438				45	1		118.4	1.95	17	
14	1.00 : 0.00 : 0.00 → 1	0.004	1.00 : 0.00 : 0.00 → 1	0.000	1.00 : 0.00 : 0.00 → 1	0.003				52	1		50.3	1.78	17	
15	0.18 : 0.82 : 0.00 → 2	0.817	0.10 : 0.90 : 0.00 → 2	0.899	0.36 : 0.64 : 0.00 → 2	0.640				36	1		99.9	1.99	18	
16	0.00 : 0.72 : 0.28 → 2	0.720	0.00 : 0.00 : 1.00 → 3	0.000	0.00 : 0.01 : 0.98 → 3	0.015				44	0		59.1	1.52	17	
17	0.19 : 0.81 : 0.00 → 2	0.811	0.27 : 0.72 : 0.01 → 2	0.731	0.34 : 0.65 : 0.01 → 2	0.659				46	1		108.8	1.64	18	
18	0.99 : 0.01 : 0.00 → 1	0.014	0.99 : 0.01 : 0.00 → 1	0.013	0.87 : 0.13 : 0.00 → 1	0.132				35	1		100.9	1.63	16	
19	0.34 : 0.66 : 0.00 → 2	0.656	0.26 : 0.74 : 0.00 → 2	0.739	0.34 : 0.66 : 0.00 → 2	0.658				19	0		65.2	1.52	18	
20	0.55 : 0.45 : 0.00 → 1	0.453	0.51 : 0.49 : 0.00 → 1	0.494	0.45 : 0.55 : 0.00 → 2	0.546				42	0		56.7	1.53	19	
Show performance scores										Target class: (Average over classes)						
										Model	AUC	CA	F1	Prec	Recall	MCC
										Logistic Regression	0.959	0.849	0.848	0.850	0.849	0.758
										Naive Bayes	0.956	0.873	0.871	0.876	0.873	0.800
										Random Forest	0.960	0.883	0.881	0.890	0.883	0.818

Conclusiones y Lecciones Aprendidas

A lo largo del proyecto aprendimos que, con un buen tratamiento de los datos, es posible predecir de forma bastante precisa el nivel de experiencia de cada usuario a partir de su comportamiento dentro del gimnasio. El modelo de Random Forest fue el que mejor captó estas diferencias, mostrando relaciones claras entre las variables de entrenamiento y el nivel de experiencia.

Los resultados confirmaron que la frecuencia de entrenamiento, la duración de las sesiones y el gasto calórico son los factores que más pesan al momento de determinar si una persona es principiante, intermedia o avanzada. También aprendimos que, aunque ciertas métricas (como la AUC o la accuracy) son útiles, en este tipo de problema importa más reducir los falsos negativos, es decir, evitar clasificar a un usuario avanzado como principiante o viceversa, porque eso impacta directamente en la calidad de las rutinas que se le asignan.

Desde una mirada más práctica, el trabajo muestra cómo este tipo de modelo podría ayudar a un gimnasio a ajustar las rutinas según el nivel real de cada persona. Por ejemplo, al identificar correctamente a los principiantes, se podrían evitar entrenamientos demasiado exigentes que generen frustración y abandono. En cambio, reconocer a los usuarios avanzados permitiría ofrecerles desafíos adecuados a su capacidad, mejorando su experiencia y motivación.