

# Predicciones de salidas de colaboradores de RiotGames (Caso Hipotetico)

## HR Employee Attrition Predicton on RiotGames Company (Hypothetical Case)

### Contexto/Context

#### ESPAÑOL

Riot Games es una multinacional con miles de colaboradores repartidos por todo el mundo. La empresa anehela la contratación de los mejores talentos disponibles y en retenerlos el mayor tiempo posible. Por lo cual, se invierte una gran cantidad de recurso (y dinero) en retener a los empleados existentes a traves de diversas iniciativas. Sin embargo, la organización (gerente de operaciones) quiere reducir los costos de retener empleados. Para ello, nos proponenn limitar los incentivos unicamente a los empleados que corren riesgo de abandono.

Se le otorgo el desafio a la area de people de identificar patrones y características de los empleados que abandonan la organización. Ademas, deben de utilizar la información levantada para predecir si un empleado esta en riesgo de abandono.

#### ENGLISH

Riot Games is a multinational company with thousands of employees around the world. The company strives to recruit the best talent available and to retain them for as long as possible. Therefore, a great deal of resources (and money) is invested in retaining existing employees through various initiatives. However, the organization (operations manager) wants to reduce the costs of retaining employees. To do so, we propose to limit incentives to only those employees who are at risk of leaving.

The people area was given the challenge of identifying patterns and characteristics of employees who leave the organization. In addition, they must use the information gathered to predict whether an employee is at risk of leaving.

### Objetivo/Objective

#### ESPAÑOL

- 1. Identificar las distintas variables o factores que hacen a los colaboradores renunciar.
- 2. Construir un modelo que sea capaz de predecir si un colaborador va a renunciar o no.

#### ENGLISH

- 1. To identify the different factors that drive attrition in riot games employees.
- 2. To build a model to predict if an employee wil attrite or not.

### Let’s goooo

```
In [102]: 1 import pandas as pd
2 import numpy as np
3 import matplotlib.pyplot as plt
4 import seaborn as sns
5
6 #Para crecer Los datos usando puntaje Z
7 from sklearn.preprocessing import StandardScaler
8 from sklearn.model_selection import train_test_split
9
10 #Algortimos que hay que utilizar
11 from sklearn.discriminant_analysis import LinearDiscriminantAnalysis
12 from sklearn.discriminant_analysis import QuadraticDiscriminantAnalysis
13 from sklearn.linear_model import LogisticRegression
14 from sklearn.neighbors import KNeighborsClassifier
15
16 #Metricas para evluar el modelo
17 from sklearn.metrics import confusion_matrix, classification_report, precision_recall_curve
18
19 #Para afinar el modelo
20 from sklearn.model_selection import GridSearchCV
21
22 #Para ignorar Las warnings
23 import warnings
24 warnings.filterwarnings("ignore")
```

### Aqui agregamos la Data / Load the data

Debes subir tu excell con la base de datos que recolectaste de tu propia organización. En este caso utilizaremos una base ficticia de Riot Games.

```
In [103]: 1 df = pd.read_excel(r'C:\Users\marti\OneDrive\Escritorio\Importante\Data Science -MIT\HR_RiotGames_Employee_Attrition_Dataset.xlsx')
```

```
In [104]: 1 df.head(10)
```

	EmployeeNumber	Attrition	Age	BusinessTravel	DailyRate	Department	DistanceFromHome	Education	EducationField	EnvironmentSatisfaction	...	StandardHours	StockOptionLevel	TotalWorkingYears	TrainingTimesLastYear	WorkLifeBala
0	1	Yes	41	Travel_Rarely	1102	Business Department	1	2	Programming	2	...	80	0	8	0	
1	2	No	49	Travel_Frequently	279	Technology Department	8	1	Programming	3	...	80	1	10	3	
2	3	Yes	37	Travel_Rarely	1373	Technology Department	2	2	Other	4	...	80	0	7	3	
3	4	No	33	Travel_Frequently	1392	Technology Department	3	4	Programming	4	...	80	0	8	3	
4	5	No	27	Travel_Rarely	591	Technology Department	2	1	Business	1	...	80	1	6	3	
5	6	No	32	Travel_Frequently	1005	Technology Department	2	2	Programming	4	...	80	0	8	2	
6	7	No	59	Travel_Rarely	1324	Technology Department	3	3	Business	3	...	80	3	12	3	
7	8	No	30	Travel_Rarely	1358	Technology Department	24	1	Programming	4	...	80	1	1	2	
8	9	No	38	Travel_Frequently	216	Technology Department	23	3	Programming	4	...	80	0	10	2	
9	10	No	36	Travel_Rarely	1299	Technology Department	27	3	Business	3	...	80	2	17	3	

10 rows × 35 columns

```
In [105]: 1 df.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 2940 entries, 0 to 2939
Data columns (total 35 columns):
#   Column                Non-Null Count  Dtype
---  -
0   EmployeeNumber        2940 non-null   int64
1   Attrition              2940 non-null   object
2   Age                    2940 non-null   int64
3   BusinessTravel         2940 non-null   object
4   DailyRate              2940 non-null   int64
5   Department             2940 non-null   object
6   DistanceFromHome       2940 non-null   int64
7   Education              2940 non-null   int64
8   EducationField         2940 non-null   object
9   EnvironmentSatisfaction 2940 non-null   int64
10  Gender                 2940 non-null   object
11  HourlyRate             2940 non-null   int64
12  JobInvolvement         2940 non-null   int64
13  JobLevel               2940 non-null   int64
14  JobRole                2940 non-null   object
15  JobSatisfaction         2940 non-null   int64
16  MaritalStatus          2940 non-null   object
17  MonthlyIncome          2940 non-null   int64
18  MonthlyRate            2940 non-null   int64
19  NumCompaniesWorked     2940 non-null   int64
20  Over18                 2940 non-null   object
21  OverTime               2940 non-null   object
22  PercentSalaryHike      2940 non-null   int64
23  PerformanceRating      2940 non-null   int64
24  RelationshipSatisfaction 2940 non-null   int64
25  StandardHours          2940 non-null   int64
26  StockOptionLevel       2940 non-null   int64
27  TotalWorkingYears      2940 non-null   int64
28  TrainingTimesLastYear  2940 non-null   int64
29  WorkLifeBalance        2940 non-null   int64
30  YearsAtCompany         2940 non-null   int64
31  YearsInCurrentRole     2940 non-null   int64
32  YearsSinceLastPromotion 2940 non-null   int64
33  YearsWithCurrManager   2940 non-null   int64
34  Gamers                 2940 non-null   object

dtypes: int64(25), object(10)
memory usage: 884.0+ KB
```

¿Que vemos? Nos permite saber que cosas son Int64 (numeros) y que cosas son Object (letras)

- 1. Hay un total de 2949 observaciones y 35 columnas
- 2. Todas las columnas son "non-null", por ende, no hay data que falte en el excel.

```
In [106]: 1 df.nunique ()

Out[106]: EmployeeNumber        2940
Attrition                2
Age                      43
BusinessTravel            3
DailyRate                886
Department                3
DistanceFromHome         29
Education                 5
EducationField            6
EnvironmentSatisfaction   4
Gender                   2
HourlyRate               71
JobInvolvement            4
JobLevel                  5
JobRole                   9
JobSatisfaction           4
MaritalStatus             3
MonthlyIncome            1349
MonthlyRate              1427
NumCompaniesWorked       10
Over18                    1
OverTime                  2
PercentSalaryHike        15
PerformanceRating         2
RelationshipSatisfaction   4
StandardHours             1
StockOptionLevel          4
TotalWorkingYears        40
TrainingTimesLastYear     7
WorkLifeBalance           4
YearsAtCompany            37
YearsInCurrentRole        19
YearsSinceLastPromotion   16
YearsWithCurrManager      18
Gamers                    2
dtype: int64
```

¿Que vemos? Aqui se muestran el numero de "opciones" que tiene cada una de las categorias.

- 1. Employee Number es el codigo de cada colaborador, es algo unico. Asi que esta columna no nos va a gregar ningun valor.
- 2. Las columnas de Over18 y StandarHours solo tienen un valor, asi que estas columnas tambien no nos agregaran valor.
- 3. Con estos numeros podemos identificar cuales columnas son continuas y cuales son categoricas.

```
In [107]: 1 #Columnas que seran eliminadas
2 df=df.drop(['EmployeeNumber','Over18','StandardHours'], axis=1)

In [108]: 1 #Creando columnas numericas

In [109]: 1 #Creando columnas numericas
2 num_cols=['DailyRate','Age','DistanceFromHome','MonthlyIncome','MonthlyRate','PercentSalaryHike','TotalWorkingYears','YearsAtCompany','NumCompaniesWorked','HourlyRate','YearsInCurrentRole','Y
3
4 #Creando columnas categoricas
5 cat_cols=['Attrition','OverTime','BusinessTravel','Department','Education','EducationField','JobSatisfaction','EnvironmentSatisfaction','WorkLifeBalance','StockOptionLevel','Gender','Performar
6
7 < >
```

## A explorar y analizar la data / Exploratory Data Analysis

ANALISIS A LAS COLUMNAS NUMERICAS Y VER QUE SACAMOS DE AHI

```
In [110]: 1 #Resumen estadístico...
2 df[num_cols].describe().T
```

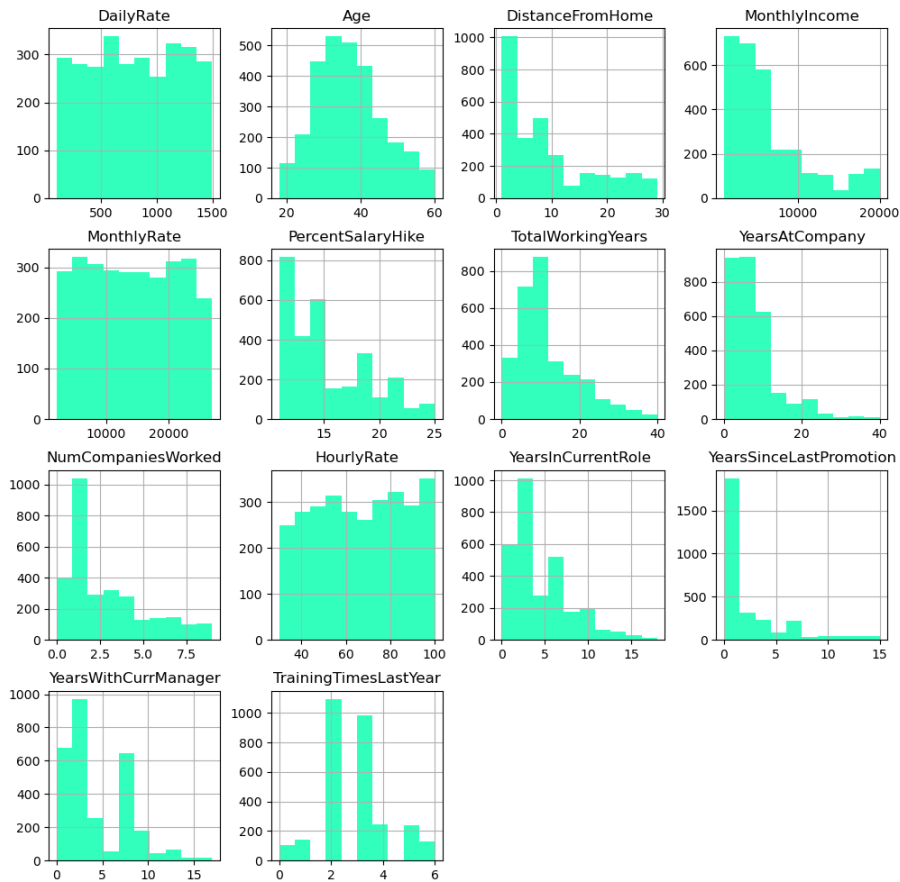
```
Out[110]:
```

	count	mean	std	min	25%	50%	75%	max
DailyRate	2940.0	802.485714	403.440447	102.0	465.0	802.0	1157.0	1499.0
Age	2940.0	36.923810	9.133819	18.0	30.0	36.0	43.0	60.0
DistanceFromHome	2940.0	9.192517	8.105485	1.0	2.0	7.0	14.0	29.0
MonthlyIncome	2940.0	6502.931293	4707.155770	1009.0	2911.0	4919.0	8380.0	19999.0
MonthlyRate	2940.0	14313.103401	7116.575021	2094.0	8045.0	14235.5	20462.0	26999.0
PercentSalaryHike	2940.0	15.209524	3.659315	11.0	12.0	14.0	18.0	25.0
TotalWorkingYears	2940.0	11.279592	7.779458	0.0	6.0	10.0	15.0	40.0
YearsAtCompany	2940.0	7.008163	6.125483	0.0	3.0	5.0	9.0	40.0
NumCompaniesWorked	2940.0	2.693197	2.497584	0.0	1.0	2.0	4.0	9.0
HourlyRate	2940.0	65.891156	20.325969	30.0	48.0	66.0	84.0	100.0
YearsInCurrentRole	2940.0	4.229252	3.622521	0.0	2.0	3.0	7.0	18.0
YearsSinceLastPromotion	2940.0	2.187755	3.221882	0.0	0.0	1.0	3.0	15.0
YearsWithCurrManager	2940.0	4.123129	3.567529	0.0	2.0	3.0	7.0	17.0
TrainingTimesLastYear	2940.0	2.799320	1.289051	0.0	2.0	3.0	3.0	6.0

#### Observaciones:

- La edad media de los empleados esta alrededor de los 37 años. Tienen una muy buena diversidad de edades en la empresa, ya que tienen un rango entre los 18 años y 60 años en los colaboradores.
- Al menos el 50% de los empleados viven en un radio de 7km de la organización (consideremos para este caso que RiotGames cuenta con solo una sede, que es su casa matriz en Los Angeles). El colaborador mas lejano esta a 29km.
- El sueldo mensual medio de un empleado de RiotGames es de 6500USD. Los sueldos varían entre 1k-20k USD, lo cual es bastante razonable. Si se ven los cuartiles y el aumento de sueldo recibido en cada uno de estos, se puede observar una desproporcion entre el tercer cuartil (75%) y el valor maximo, donde el primero gana en promedio 8.400USD y el maximo gana un promedio de 20.000USD. En resumen, **los que mas ganan en la empresa tienen ingresos desproporcionados con el resto (lo cual es bastante comun en las organización).**
- La subida salarial media de un colaborador ronda en el 15% y al menos el 50% de los colaboradores obtuvo un aumento salarial del 14% o menos. Me quito el sombrero, demuestra que tienen un buen modelo de compensaciones en la organización y que se esta pendiente de esto.
- El tiempo transcurrido desde que un colaborador obtuvo un ascenso es de ~2,19. La mayoría de los empleados han sido ascendidos en el último año.
- La media de años vinculado en la empresa los colaboradores es de 7.

```
In [111]: 1 #Ahora veremos la distribución de los datos. Crearemos histogramas
2 df[num_cols].hist(figsize=(12,12), color='#33FF80')
3 plt.show()
```



#### Observaciones:

- La distribución de la edad de los colaboradores se parece a una distribución normal. Donde la mayoría esta entre los 25 y 50.
- La mayoría de los colaboradores viven cerca de la casa matriz y cada vez va bajando mas la gente en relacion a la distancia.
- Los ingresos mensuales y el número total de años de trabajo están sesgados a la derecha, lo que indica que la mayoría de los trabajadores ocupan puestos de nivel inicial o medio en la empresa.
- El porcentaje de aumento salarial está desviado hacia la derecha, lo que significa que la mayoría de los empleados reciben un porcentaje de aumento salarial más bajo.
- Existe una buena proporción en YearsAtCompany respecto a los colaboradores que llevan mas de 10 años. Lo que demuestra que hay varias personas leales a la compañía.
- La distribución de la variable YearsSinceLastPromotion indica que algunos empleados no han recibido un ascenso en 10-15 años y siguen trabajando en la empresa. Se supone que estos empleados tienen una gran experiencia laboral y ocupan puestos de alta dirección, como cofundadores, empleados de la C-suite, etc (pero abría que indagar estos perfiles).

#### ANALISIS A LAS COLUMNAS CATEGORICAS Y VER QUE SACAMOS DE AHÍ

Con esto podremos ver el porcentaje de las categorías y cuantas existen por "opcion".

```
In [112]: 1 for i in cat_cols:
2         print(df[i].value_counts(normalize = True))
3         print('*'*40)
```

```
No      0.838776
Yes      0.161224
Name: Attrition, dtype: float64
*****
No      0.717007
Yes      0.282993
Name: OverTime, dtype: float64
*****
Travel_Rarely      0.709524
Travel_Frequently  0.188435
Non-Travel         0.102041
Name: BusinessTravel, dtype: float64
*****
Technology Department      0.653741
Business Department        0.303401
Human Resources             0.042857
Name: Department, dtype: float64
*****
3      0.389116
4      0.270748
2      0.191837
1      0.115646
5      0.032653
Name: Education, dtype: float64
*****
Programming      0.412245
Business          0.315646
Marketing         0.108163
Technical Degree  0.089796
Other             0.055782
Human Resources   0.018367
Name: EducationField, dtype: float64
*****
4      0.312245
3      0.300680
1      0.196599
2      0.190476
Name: JobSatisfaction, dtype: float64
*****
3      0.308163
4      0.303401
2      0.195238
1      0.193197
Name: EnvironmentSatisfaction, dtype: float64
*****
3      0.607483
2      0.234014
4      0.104082
1      0.054422
Name: WorkLifeBalance, dtype: float64
*****
0      0.429252
1      0.405442
2      0.107483
3      0.057823
Name: StockOptionLevel, dtype: float64
*****
Male      0.6
Female    0.4
Name: Gender, dtype: float64
*****
3      0.846259
4      0.153741
Name: PerformanceRating, dtype: float64
*****
3      0.590476
2      0.255102
4      0.097959
1      0.056463
Name: JobInvolvement, dtype: float64
*****
1      0.369388
2      0.363265
3      0.148299
4      0.072109
5      0.046939
Name: JobLevel, dtype: float64
*****
Marketing Analyst      0.221769
Designer              0.198639
Developer             0.176190
Manufacturing Director 0.098639
Customer Experience    0.089116
Manager              0.069388
Business Analyst       0.056463
Research Director     0.054422
Human Resources        0.035374
Name: JobRole, dtype: float64
*****
Married      0.457823
Single       0.319728
Divorced     0.222449
Name: MaritalStatus, dtype: float64
*****
3      0.312245
4      0.293878
2      0.206122
1      0.187755
Name: RelationshipSatisfaction, dtype: float64
*****
Yes      0.727551
No       0.272449
Name: Gamers, dtype: float64
*****
```

Observaciones:

- La tasa de bajas de colaboradores es del 16%.
- Alrededor de 28% de los trabajadores trabajan horas extras. Un gran numero que influye en el balance de calidad de vida personal-profesional (factor de estres).
- Aproximadamente un 73% de los colaboradores vienen del mundo de la programación y business. Y el 65% de los colaboradores trabajan en el departamento de tecnología.
- Cerca del 40% de los trabajadores evalúan su satisfacción laboral o ambiente laboral como baja o media. Lo cual es un gran porcentaje de la plantilla.
- 19% de los colaboradores viaja frecuentemente por el trabajo y 71% de los trabajadores rara vez.
- Mas del 30% de los empleados muestran una implicancia baja o media. Y mas del 80% de los empleados no tiene stockoptions o tiene muy pocas.
- El 73% de los colaboradores de RiotGames tiene relación con el mundo gamers. Lo cual podría implicar relación con el producto y cultura de la organización.
- Ninguno de los empleados a recibido una puntuación menor a 3 (excelente) respecto a su evaluación de desempeño (rendimiento). El 85% dfe sus colaboradores es evaluado 3 (excelente) y el resto 4 (sobre lo esperado). Dos escenarios:
  - Los colaboradores de Riot Games son increíblemente buenos.
  - El proceso de evaluación es indulgente. Y falta capacitar a la organizacion y lideres de como generar una correcta evaluación de desempeño.

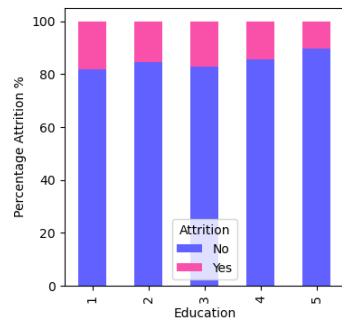
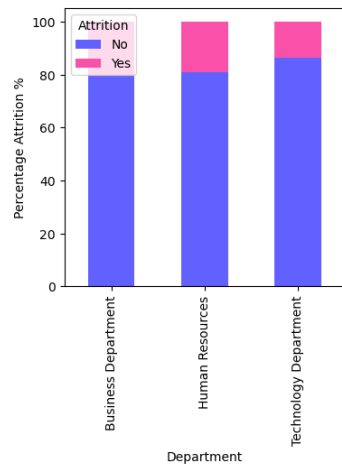
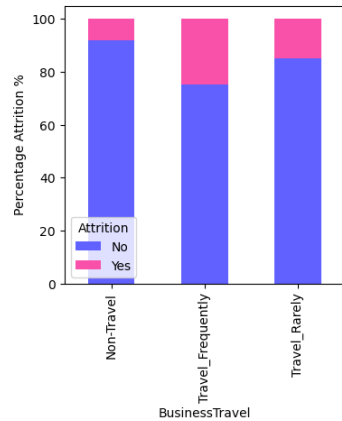
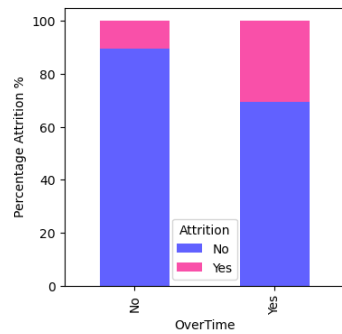
## MixUp (Análisis columnas categóricas) / Multivariate Analysis

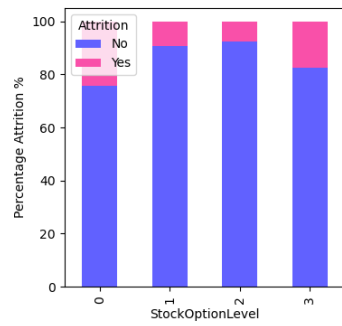
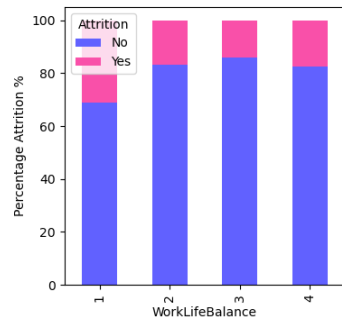
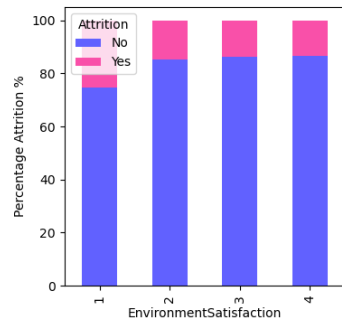
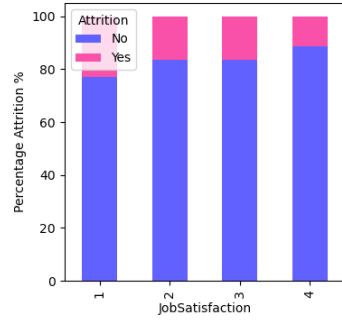
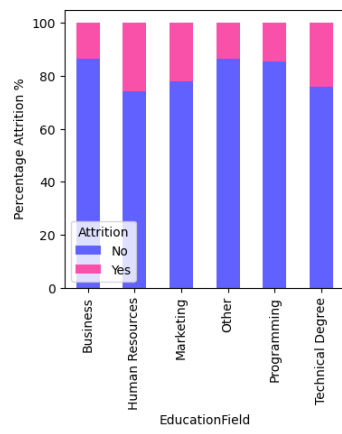
Ahora analizaremos el relación entre la tasa de abandono (attrition) con todas las otras variables categóricas.

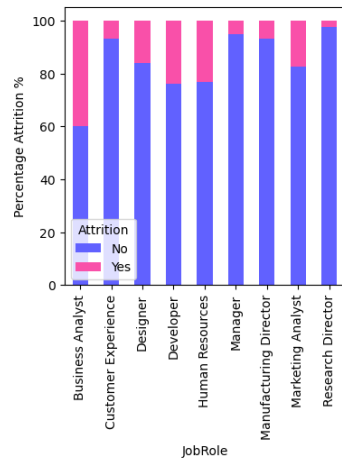
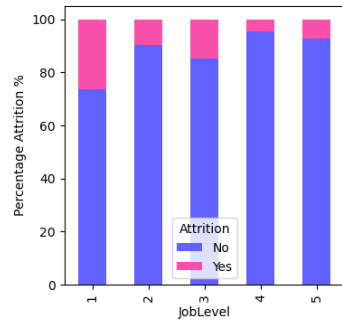
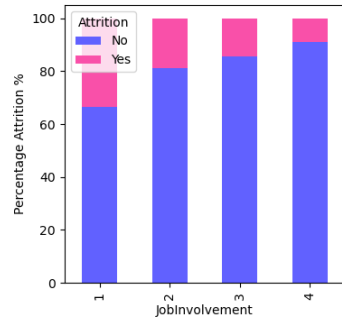
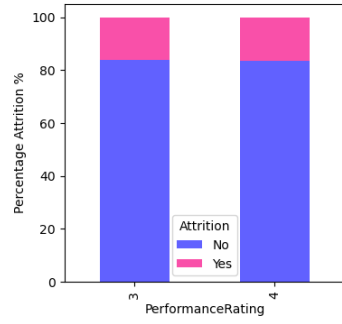
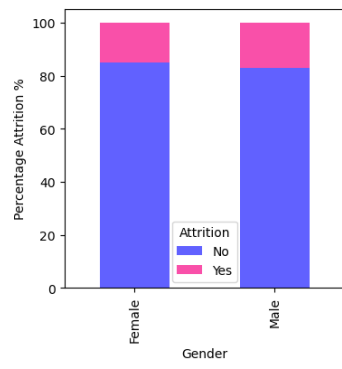
```

In [113]: 1 for i in cat_cols:
2         if i != 'Attrition':
3             (pd.crosstab(df[i], df['Attrition'], normalize = 'index')*100).plot(kind='bar', figsize=(4,4), stacked= True, color = ['#6161FF', '#F958A9'])
4             plt.ylabel('Percentage Attrition %')

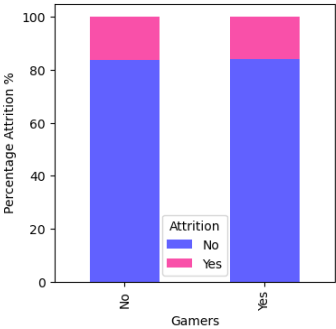
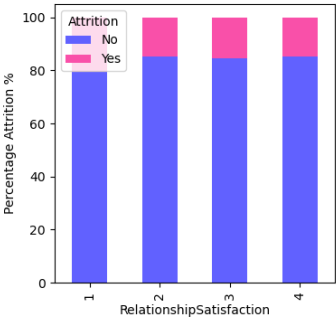
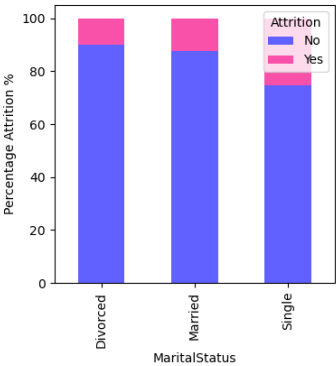
```











Observaciones:

- Los trabajadores que hacen horas extra tienen un 30% más de probabilidad de abandono. Lo cual es alto si se compara con el 10% de las personas que no tienen overtime.
- Como se comentó antes, la mayoría trabaja en el departamento de tecnología y es este mismo departamento con el menor índice de rotación. La probabilidad de que se produzcan bajas es del ~15%.
- Los cargos HR, Technical degree y Marketin son los que muestran un mayor índice de rotación del personal. Se podría hipotetizar que se prioriza la retención de cargos que generan mas ingreso a la compañía y se descuidan los cargos que generan un mayor "costo" a la organización. Ya que los devs y los de business son los que menos rotación tienen y a la vez son esenciales para desarrollar y crear los productos.
- Se demuestra un increíble correlación entre el grado de implicancia y el nivel de bajas de colaboradores. Donde **entre menos implicancia mayor es el porcentaje de rotación**. Donde los que estan en el nivel mas bajo de implicancia tienen un porcentaje del 35% de abandono. Seria interesante evaluar la relación entre implicancia y stock options, ver si existe una relación entre estos factores y el nivel de rotación. Ya que los que no tienen stock options son los que mas rotan, pero despues les continua los que más tienen stock options ¿A quien les estamos dando las stock options? ¿Tiene relación con el grado de implicancia? ¿Tener más StockOptions significa tener más implicancia con el trabajo?
- Los empleados con un nivel de empleo inferior (entry level) también sufren más bajas, y los empleados con un nivel de empleo de 1 tienen casi un 25% de probabilidades de abandonar. Puede tratarse de empleados jóvenes que tienden a explorar más opciones en las etapas iniciales de sus carreras.
- Que sean del mundo gamers no demuestra implicancia en el nivel de rotación del personal. No es un factor evidentemente fuerte en el nivel de retencion.

MixUp (Análisis columnas numericas con Attrition)

Vamos a ver la relación entre la rotación y cada una de las variables numericas!

In [114]: 1 df.groupby(['Attrition'])[num\_cols].mean()

Out[114]:

Attrition	DailyRate	Age	DistanceFromHome	MonthlyIncome	MonthlyRate	PercentSalaryHike	TotalWorkingYears	YearsAtCompany	NumCompaniesWorked	HourlyRate	YearsInCurrentRole	YearsSinceLastPromotion	YearsWithCurrManage
No	812.504461	37.561233	8.915653	6832.739659	14265.779400	15.231144	11.862936	7.369019	2.645580	65.952149	4.484185	2.234388	4.367391
Yes	750.362869	33.607595	10.632911	4787.092827	14559.308017	15.097046	8.244726	5.130802	2.940928	65.573840	2.902954	1.945148	2.852321

Observaciones:

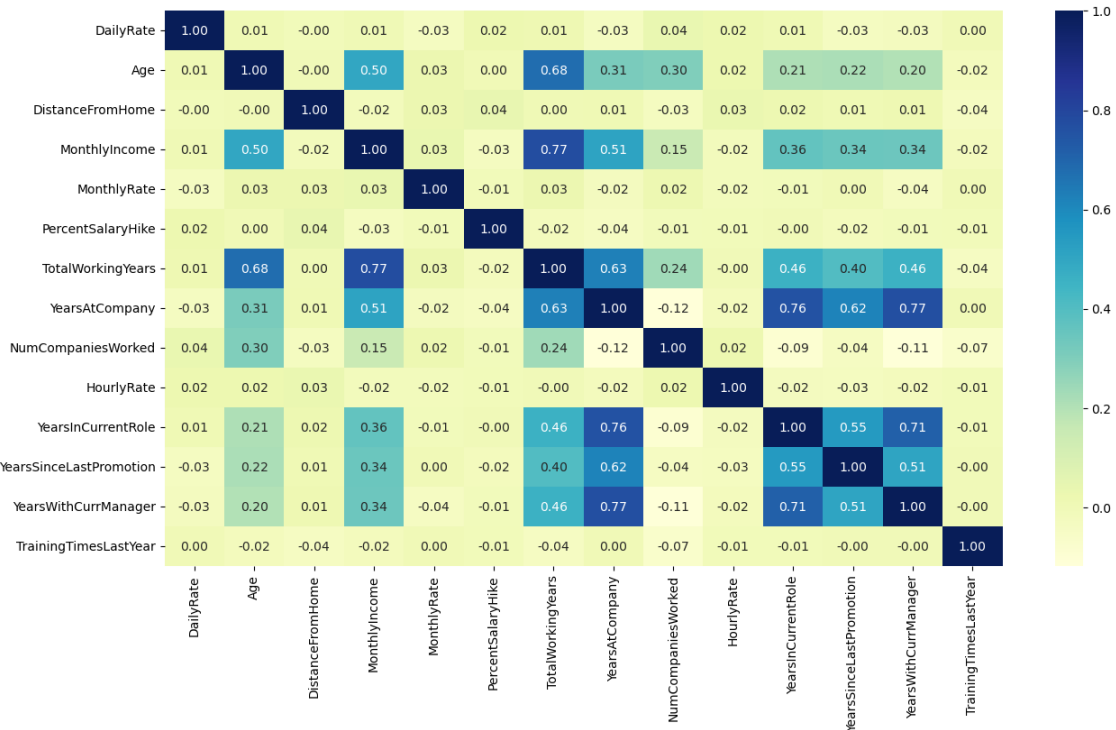
- Los colaboradores que abandonan la empresa tienen casi un 30% menos de ingresos medios y un 30% menos de experiencia laboral de los que no lo hacen. Se puede hipotetizar que lo hacen para explorar nuevas opciones y/o aumentar su salario con un cambio de empresa.
- La distancia es relevante? Los empleados que se dan de baja también tienden a vivir un 16% más lejos de la oficina que los que no se dan de baja. ¿Si se aplica teletrabajo como se vera afectado esto?

ANALISIS DE RELACIÓN ENTRE LAS DISTINTAS VARIABLES NUMERICAS

Vamos a ver como se correlacionan entre ellas. Donde haremos un mapa de calor y tendremos que identificar relaciones interesantes que puedan significar algo.

```
In [115]: 1 plt.figure(figsize=(15,8))
2
3 sns.heatmap(df[num_cols].corr(), annot = True, fmt = '0.2f', cmap = 'YlGnBu')

Out[115]: <Axes: >
```



Observaciones:

- En las correlación es importante ver que relación de variables se acerca al 1.0 (entre mas cercano a ese numero mayor es la correlación).
- Experiencia laboral, sueldo mensual, años en Riot Games y años con el mismo líder estan altamente correlacionado entre si y con la edad de los empleados**, en donde el último es bastante logico, ya que a medida que uno va teniendo mas años tiende a tener mas de las otras variables.
- Los años en la empresa y los años en el puesto actual están correlacionados con los años transcurridos desde la última promoción**, lo que significa que la empresa no está dando promociones en el momento adecuado. Es decir, entre mas años llevo en el cargo actual, más tiempo ha pasado desde mi ultima promoción.

Resumen del analsis de datos

- La edad media de los empleados ronda los 37 años. Tiene un rango elevado, de 18 años a 60, lo que indica una buena diversidad de edades en la organización.
- Al menos el 50% de los empleados viven en un radio de 7 km de la organización. Sin embargo, hay algunos valores extremos, ya que el valor máximo es de 29 km.
- El ingreso mensual medio de un empleado es de 6500 USD. Tiene un alto rango de valores de 1K-20K USD, lo que es de esperar para la distribución de ingresos de cualquier organización. Hay una gran diferencia entre el valor del tercer cuartil (alrededor de 8400 USD) y el valor máximo (casi 20000 USD), lo que demuestra que las personas que más ganan en la empresa tienen unos ingresos desproporcionadamente grandes en comparación con el resto de los empleados. De nuevo, esto es bastante común en la mayoría de las organizaciones. Sin embargo, se podría optar por algunas políticas o analisis del modelo de compensación como una buena practica.
- La subida salarial media de un empleado ronda el 15%. Al menos el 50% de los empleados obtuvo un aumento salarial del 14% o menos, siendo el aumento salarial máximo del 25%.
- El número medio de años que un empleado lleva vinculado a la empresa es de 7.
- Por término medio, el número de años transcurridos desde que un empleado obtuvo un ascenso es de ~2,19. La mayoría de los empleados han sido ascendidos en el último año.
- La distribución por edades se aproxima a una distribución normal, con la mayoría de los empleados entre 25 y 50 años.
- DistanceFromHome también tiene una distribución sesgada a la derecha, lo que significa que la mayoría de los empleados viven cerca del trabajo, pero hay algunos que viven más lejos.
- Los ingresos mensuales y el número total de años de trabajo están sesgados a la derecha, lo que indica que la mayoría de los trabajadores ocupan puestos de nivel inicial o medio en la empresa.
- El porcentaje de aumento salarial está sesgado a la derecha, lo que significa que la mayoría de los empleados están recibiendo aumentos salariales de menor porcentaje.
- La distribución de la variable YearsAtCompany muestra una buena proporción de trabajadores con más de 10 años, lo que indica un número significativo de empleados leales a la organización.
- La distribución de YearsInCurrentRole tiene tres picos en 0, 2 y 7. Hay pocos empleados que hayan permanecido en el mismo puesto durante 15 años o más.
- La distribución de la variable YearsSinceLastPromotion indica que algunos empleados no han recibido un ascenso en 10-15 años y siguen trabajando en la empresa. Se supone que estos empleados tienen una gran experiencia laboral y ocupan puestos de alta dirección, como cofundadores, empleados de la C-suite, etc.
- Las distribuciones de DailyRate, HourlyRate y MonthlyRate parecen uniformes y no aportan mucha información. Podría ser que la tasa diaria se refiera a los ingresos obtenidos por día extra trabajado, mientras que la tasa horaria podría referirse al mismo concepto aplicado a las horas extra trabajadas al día. Dado que estas tasas tienden a ser muy similares para varios empleados de un mismo departamento, eso explica la distribución uniforme que muestran.
- La tasa de abandono de los empleados es del 16%.
- Alrededor del 28% de los empleados hacen horas extraordinarias. Esta cifra parece estar en el lado más alto y podría indicar una vida laboral estresada de los empleados. Los empleados que hacen horas extras tienen más de un 30% de probabilidades de abandono, una cifra muy alta comparada con el 10% de probabilidades de abandono de los empleados que no hacen horas extras.
- El 71% de los empleados ha viajado pocas veces, mientras que alrededor del 19% tiene que viajar con frecuencia.
- Alrededor del 73% de los empleados tienen formación en el campo de las Programming y Business.
- Más del 65% de los empleados trabajan en el departamento de Technology Department y Business Department. La probabilidad de desgaste es de ~15%.
- Casi el 40% de los empleados tienen una satisfacción baja (1) o media (2) con el trabajo y el entorno de la organización, lo que indica que la moral de la empresa parece ser algo baja.
- Más del 30% de los empleados muestran una implicación en el trabajo de baja (1) a media (2).
- Más del 80% de los empleados no tienen opciones sobre acciones o tienen muy pocas.
- En cuanto a la valoración del rendimiento, ninguno de los empleados ha recibido una valoración inferior a 3 (excelente). Alrededor del 85% de los empleados tiene una valoración del rendimiento igual a 3 (excelente), mientras que el resto tiene una valoración de 4 (sobresaliente). Esto podría significar que la mayoría de los empleados son de alto rendimiento, o lo más probable es que la organización sea muy indulgente con su proceso de evaluación del rendimiento. Recomendación de capacitación a líderes respecto a la nivelación de puntuación.
- Los empleados que trabajan como Business Analyst tienen una tasa de abandono de alrededor del 40%, mientras que los de RR.HH. y los Developer tienen una tasa de abandono de alrededor del 25%. Los departamentos de business tiene tasa de toración mas alta que los cargos de tecnología, lo cual es bastante raro en este nivel de industria. Habria que analizar en profundidad la area de Businees (remuneración, liderazgo, distribución de proyectos, etc.).
- Cuanto menor es la implicación en el trabajo del empleado, mayores parecen ser sus posibilidades de abandono, con un 35% de abandono entre los empleados con una implicación en el trabajo de 1 punto. Esto podría deberse a que los empleados con una menor implicación en el trabajo podrían sentirse excluidos o menos valorados y ya han empezado a explorar nuevas opciones, lo que conduce a una mayor tasa de abandono.
- Los empleados con un nivel de empleo inferior también sufren más bajas, y los empleados con un nivel de empleo de 1 tienen casi un 25% de probabilidades de abandonar. Puede tratarse de empleados jóvenes que tienden a explorar más opciones en las etapas iniciales de sus carreras.
- Una valoración baja de la conciliación de la vida laboral y personal lleva a los empleados a abandonar la empresa; ~30% de los que se encuentran en la categoría de valoración 1 presentan bajas.
- Los empleados que abandonan la empresa tienen una media de ingresos casi un 30% inferior y un 30% menos de experiencia laboral que los que no lo hacen. Estos podrían ser los empleados que buscan explorar nuevas opciones y/o aumentar su salario con un cambio de empresa.
- Los empleados que se dan de baja también tienden a vivir un 16% más lejos de la oficina que los que no se dan de baja. Los desplazamientos más largos para ir y volver del trabajo podrían significar que tienen que dedicar más tiempo/dinero cada día, y esto podría estar provocando insatisfacción laboral y el deseo de abandonar la organización.
- La experiencia laboral total, los ingresos mensuales, los años en la empresa y los años con los jefes actuales están muy correlacionados entre sí y con la edad del empleado, lo cual es fácil de entender, ya que estas variables muestran un aumento con la edad para la mayoría de los empleados.
- Los años en la empresa y los años en el puesto actual están correlacionados con los años transcurridos desde el último ascenso, lo que significa que la empresa no concede los ascensos en el momento adecuado.

TERMINA LA PARTE I (ANALISIS DE LOS DATOS) Y COMIENZA PARTE II (ALGORITMO PARA PREDECIR SALIDAS)

Modelo de Predicción de Salidas / Model Building - Approach

## HEMOS ANALIZADO LA DATA, AHORA VAMOS A CONTRUIR EL MODELO PARA PREDECIR LAS SALIDAS DE LOS COLABORADORES DE RIOTGAMES

De aqui en adelante el contenido es mas complicado de explicar, pero hare el mejor esfuerzo.

1. Preparar la dara para modelar.
2. Dividir los datos en conjuntos de entrenamientos y de test.
3. Construir el modelo con los datos de entrenamiento.
4. Ajustar el modelo.
5. Probar los datos en el conjunto de test.

## Preparando los datos para el modelo

Crando "dummy" variables para las variables categoricas. En general, es cambiar palabras por numeros.

```
In [116]: 1 # Creando La lista de columna en donde debemos crear dummy variables.
2 to_get_dummies_for = ['BusinessTravel', 'Department', 'Education', 'EducationField', 'EnvironmentSatisfaction', 'Gender', 'JobInvolvement', 'JobLevel', 'JobRole', 'MaritalStatus', 'Gamers']
3
4 # Creando Las dummy variables
5 df = pd.get_dummies(data = df, columns = to_get_dummies_for, drop_first = True)
6
7 # Mapeando overtime and attrition
8 dict_OverTime = {'Yes': 1, 'No': 0}
9 dict_attrition = {'Yes': 1, 'No': 0}
10
11 df['OverTime'] = df.OverTime.map(dict_OverTime)
12 df['Attrition'] = df.Attrition.map(dict_attrition)
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65
66
67
68
69
70
71
72
73
74
75
76
77
78
79
80
81
82
83
84
85
86
87
88
89
90
91
92
93
94
95
96
97
98
99
100
101
102
103
104
105
106
107
108
109
110
111
112
113
114
115
116
117
118
119
120
121
122
123
124
125
126
127
128
129
130
131
132
133
134
135
136
137
138
139
140
141
142
143
144
145
146
147
148
149
150
151
152
153
154
155
156
157
158
159
160
161
162
163
164
165
166
167
168
169
170
171
172
173
174
175
176
177
178
179
180
181
182
183
184
185
186
187
188
189
190
191
192
193
194
195
196
197
198
199
200
201
202
203
204
205
206
207
208
209
210
211
212
213
214
215
216
217
218
219
220
221
222
223
224
225
226
227
228
229
230
231
232
233
234
235
236
237
238
239
240
241
242
243
244
245
246
247
248
249
250
251
252
253
254
255
256
257
258
259
260
261
262
263
264
265
266
267
268
269
270
271
272
273
274
275
276
277
278
279
280
281
282
283
284
285
286
287
288
289
290
291
292
293
294
295
296
297
298
299
300
301
302
303
304
305
306
307
308
309
310
311
312
313
314
315
316
317
318
319
320
321
322
323
324
325
326
327
328
329
330
331
332
333
334
335
336
337
338
339
340
341
342
343
344
345
346
347
348
349
350
351
352
353
354
355
356
357
358
359
360
361
362
363
364
365
366
367
368
369
370
371
372
373
374
375
376
377
378
379
380
381
382
383
384
385
386
387
388
389
390
391
392
393
394
395
396
397
398
399
400
401
402
403
404
405
406
407
408
409
410
411
412
413
414
415
416
417
418
419
420
421
422
423
424
425
426
427
428
429
430
431
432
433
434
435
436
437
438
439
440
441
442
443
444
445
446
447
448
449
450
451
452
453
454
455
456
457
458
459
460
461
462
463
464
465
466
467
468
469
470
471
472
473
474
475
476
477
478
479
480
481
482
483
484
485
486
487
488
489
490
491
492
493
494
495
496
497
498
499
500
501
502
503
504
505
506
507
508
509
510
511
512
513
514
515
516
517
518
519
520
521
522
523
524
525
526
527
528
529
530
531
532
533
534
535
536
537
538
539
540
541
542
543
544
545
546
547
548
549
550
551
552
553
554
555
556
557
558
559
560
561
562
563
564
565
566
567
568
569
570
571
572
573
574
575
576
577
578
579
580
581
582
583
584
585
586
587
588
589
590
591
592
593
594
595
596
597
598
599
600
601
602
603
604
605
606
607
608
609
610
611
612
613
614
615
616
617
618
619
620
621
622
623
624
625
626
627
628
629
630
631
632
633
634
635
636
637
638
639
640
641
642
643
644
645
646
647
648
649
650
651
652
653
654
655
656
657
658
659
660
661
662
663
664
665
666
667
668
669
670
671
672
673
674
675
676
677
678
679
680
681
682
683
684
685
686
687
688
689
690
691
692
693
694
695
696
697
698
699
700
701
702
703
704
705
706
707
708
709
710
711
712
713
714
715
716
717
718
719
720
721
722
723
724
725
726
727
728
729
730
731
732
733
734
735
736
737
738
739
740
741
742
743
744
745
746
747
748
749
750
751
752
753
754
755
756
757
758
759
760
761
762
763
764
765
766
767
768
769
770
771
772
773
774
775
776
777
778
779
780
781
782
783
784
785
786
787
788
789
790
791
792
793
794
795
796
797
798
799
800
801
802
803
804
805
806
807
808
809
810
811
812
813
814
815
816
817
818
819
820
821
822
823
824
825
826
827
828
829
830
831
832
833
834
835
836
837
838
839
840
841
842
843
844
845
846
847
848
849
850
851
852
853
854
855
856
857
858
859
860
861
862
863
864
865
866
867
868
869
870
871
872
873
874
875
876
877
878
879
880
881
882
883
884
885
886
887
888
889
890
891
892
893
894
895
896
897
898
899
900
901
902
903
904
905
906
907
908
909
910
911
912
913
914
915
916
917
918
919
920
921
922
923
924
925
926
927
928
929
930
931
932
933
934
935
936
937
938
939
940
941
942
943
944
945
946
947
948
949
950
951
952
953
954
955
956
957
958
959
960
961
962
963
964
965
966
967
968
969
970
971
972
973
974
975
976
977
978
979
980
981
982
983
984
985
986
987
988
989
990
991
992
993
994
995
996
997
998
999
1000
1001
1002
1003
1004
1005
1006
1007
1008
1009
1010
1011
1012
1013
1014
1015
1016
1017
1018
1019
1020
1021
1022
1023
1024
1025
1026
1027
1028
1029
1030
1031
1032
1033
1034
1035
1036
1037
1038
1039
1040
1041
1042
1043
1044
1045
1046
1047
1048
1049
1050
1051
1052
1053
1054
1055
1056
1057
1058
1059
1060
1061
1062
1063
1064
1065
1066
1067
1068
1069
1070
1071
1072
1073
1074
1075
1076
1077
1078
1079
1080
1081
1082
1083
1084
1085
1086
1087
1088
1089
1090
1091
1092
1093
1094
1095
1096
1097
1098
1099
1100
1101
1102
1103
1104
1105
1106
1107
1108
1109
1110
1111
1112
1113
1114
1115
1116
1117
1118
1119
1120
1121
1122
1123
1124
1125
1126
1127
1128
1129
1130
1131
1132
1133
1134
1135
1136
1137
1138
1139
1140
1141
1142
1143
1144
1145
1146
1147
1148
1149
1150
1151
1152
1153
1154
1155
1156
1157
1158
1159
1160
1161
1162
1163
1164
1165
1166
1167
1168
1169
1170
1171
1172
1173
1174
1175
1176
1177
1178
1179
1180
1181
1182
1183
1184
1185
1186
1187
1188
1189
1190
1191
1192
1193
1194
1195
1196
1197
1198
1199
1200
1201
1202
1203
1204
1205
1206
1207
1208
1209
1210
1211
1212
1213
1214
1215
1216
1217
1218
1219
1220
1221
1222
1223
1224
1225
1226
1227
1228
1229
1230
1231
1232
1233
1234
1235
1236
1237
1238
1239
1240
1241
1242
1243
1244
1245
1246
1247
1248
1249
1250
1251
1252
1253
1254
1255
1256
1257
1258
1259
1260
1261
1262
1263
1264
1265
1266
1267
1268
1269
1270
1271
1272
1273
1274
1275
1276
1277
1278
1279
1280
1281
1282
1283
1284
1285
1286
1287
1288
1289
1290
1291
1292
1293
1294
1295
1296
1297
1298
1299
1300
1301
1302
1303
1304
1305
1306
1307
1308
1309
1310
1311
1312
1313
1314
1315
1316
1317
1318
1319
1320
1321
1322
1323
1324
1325
1326
1327
1328
1329
1330
1331
1332
1333
1334
1335
1336
1337
1338
1339
1340
1341
1342
1343
1344
1345
1346
1347
1348
1349
1350
1351
1352
1353
1354
1355
1356
1357
1358
1359
1360
1361
1362
1363
1364
1365
1366
1367
1368
1369
1370
1371
1372
1373
1374
1375
1376
1377
1378
1379
1380
1381
1382
1383
1384
1385
1386
1387
1388
1389
1390
1391
1392
1393
1394
1395
1396
1397
1398
1399
1400
1401
1402
1403
1404
1405
1406
1407
1408
1409
1410
1411
1412
1413
1414
1415
1416
1417
1418
1419
1420
1421
1422
1423
1424
1425
1426
1427
1428
1429
1430
1431
1432
1433
1434
1435
1436
1437
1438
1439
1440
1441
1442
1443
1444
1445
1446
1447
1448
1449
1450
1451
1452
1453
1454
1455
1456
1457
1458
1459
1460
1461
1462
1463
1464
1465
1466
1467
1468
1469
1470
1471
1472
1473
1474
1475
1476
1477
1478
1479
1480
1481
1482
1483
1484
1485
1486
1487
1488
1489
1490
1491
1492
1493
1494
1495
1496
1497
1498
1499
1500
1501
1502
1503
1504
1505
1506
1507
1508
1509
1510
1511
1512
1513
1514
1515
1516
1517
1518
1519
1520
1521
1522
1523
1524
1525
1526
1527
1528
1529
1530
1531
1532
1533
1534
1535
1536
1537
1538
1539
1540
1541
1542
1543
1544
1545
1546
1547
1548
1549
1550
1551
1552
1553
1554
1555
1556
1557
1558
1559
1560
1561
1562
1563
1564
1565
1566
1567
1568
1569
1570
1571
1572
1573
1574
1575
1576
1577
1578
1579
1580
1581
1582
1583
1584
1585
1586
1587
1588
1589
1590
1591
1592
1593
1594
1595
1596
1597
1598
1599
1600
1601
1602
1603
1604
1605
1606
1607
1608
1609
1610
1611
1612
1613
1614
1615
1616
1617
1618
1619
1620
1621
1622
1623
1624
1625
1626
1627
1628
1629
1630
1631
1632
1633
1634
1635
1636
1637
1638
1639
1640
1641
1642
1643
1644
1645
1646
1647
1648
1649
1650
1651
1652
1653
1654
1655
1656
1657
1658
1659
1660
1661
1662
1663
1664
1665
1666
1667
1668
1669
1670
1671
1672
1673
1674
1675
1676
1677
1678
1679
1680
1681
1682
1683
1684
1685
1686
1687
1688
1689
1690
1691
1692
1693
1694
1695
1696
1697
1698
1699
1700
1701
1702
1703
1704
1705
1706
1707
1708
1709
1710
1711
1712
1713
1714
1715
1716
1717
1718
1719
1720
1721
1722
1723
1724
1725
1726
1727
1728
1729
1730
1731
1732
1733
1734
1735
1736
1737
1738
1739
1740
1741
1742
1743
1744
1745
1746
1747
1748
1749
1750
1751
1752
1753
1754
1755
1756
1757
1758
1759
1760
1761
1762
1763
1764
1765
1766
1767
1768
1769
1770
1771
1772
1773
1774
1775
1776
1777
1778
1779
1780
1781
1782
1783
1784
1785
1786
1787
1788
1789
1790
1791
1792
1793
1794
1795
1796
1797
1798
1799
1800
1801
1802
1803
1804
1805
1806
1807
1808
1809
1810
1811
1812
1813
1814
1815
1816
1817
1818
1819
1820
1821
1822
1823
1824
1825
1826
1827
1828
1829
1830
1831
1832
1833
1834
1835
1836
1837
1838
1839
1840
1841
1842
1843
1844
1845
1846
1847
1848
1849
1850
1851
1852
1853
1854
1855
1856
1857
1858
1859
1860
1861
1862
1863
1864
1865
1866
1867
1868
1869
1870
1871
1872
1873
1874
1875
1876
1877
1878
1879
1880
1881
1882
1883
1884
1885
1886
1887
1888
1889
1890
1891
1892
1893
1894
1895
1896
1897
1898
1899
1900
1901
1902
1903
1904
1905
1906
1907
1908
1909
1910
1911
1912
1913
1914
1915
1916
1917
1918
1919
1920
1921
1922
1923
1924
1925
1926
1927
1928
1929
1930
1931
1932
1933
1934
1935
1936
1937
1938
1939
1940
1941
1942
1943
1944
1945
1946
1947
1948
1949
1950
1951
1952
1953
1954
1955
1956
1957
1958
1959
1960
1961
1962
1963
1964
1965
1966
1967
1968
1969
1970
1971
1972
1973
1974
1975
1976
1977
1978
1979
1980
1981
1982
1983
1984
1985
1986
1987
1988
1989
1990
1991
1992
1993
1994
1995
1996
1997
1998
1999
2000
2001
2002
2003
2004
2005
2006
2007
2008
2009
2010
2011
2012
2013
2014
2015
2016
2017
2018
2019
2020
2021
2022
2023
2024
2025
2026
2027
2028
2029
2030
2031
2032
2033
2034
2035
2036
2037
2038
2039
2040
2041
2042
2043
2044
2045
2046
2047
2048
2049
2050
2051
2052
2053
2054
2055
2056
2057
2058
2059
2060
2061
2062
2063
2064
2065
2066
2067
2068
2069
2070
2071
2072
2073
2074
2075
2076
2077
2078
2079
2080
2081
2082
2083
2084
2085
2086
2087
2088
2089
2090
2091
2092
2093
2094
2095
2096
2097
2098
2099
2100
2101
2102
2103
2104
2105
2106
2107
2108
2109
2110
2111
2112
2113
2114
2115
2116
2117
2118
2119
2120
2121
2122
2123
2124
2125
2126
2127
2128
2129
2130
2131
2132
2133
2134
2135
2136
2137
2138
2139
2140
2141
2142
2143
2144
2145
2146
2147
2148
2149
2150
2151
2152
2153
2154
2155
2156
2157
2158
2159
2160
2161
2162
2163
2164
2165
2166
2167
2168
2169
2170
2171
2172
2173
2174
2175
2176
2177
2178
2179
2180
2181
2182
2183
2184
2185
2186
2187
2188
2189
2190
2191
2192
2193
2194
2195
2196
2197
2198
2199
2200
2201
2202
2203
2204
2205
2206
2207
2208
2209
2210
2211
2212
2213
2214
2215
2216
2217
2218
2219
2220
2221
2222
2223
2224
2225
2226
2227
2228
2229
2230
2231
2232
2233
2234
2235
2236
2237
2238
2239
2240
2241
2242
2243
2244
2245
2246
2247
2248
2249
2250
2251
2252
2253
2254
2255
2256
2257
2258
2259
2260
2261
2262
2263
2264
2265
2266
2267
2268
2269
2270
2271
2272
2273
2274
2275
2276
2277
2278
2279
2280
2281
2282
2283
2284
2285
2286
2287
2288
2289
2290
2291
2292
2293
2294
2295
2296
2297
2298
2299
2300
2301
2302
2303
2304
2305
2306
2307
2308
2309
2310
2311
2312
2313
2314
2315
2316
2317
2318
2319
2320
2321
2322
2323
2324
2325
2326
2327
2328
2329
2330
2331
2332
2333
2334
2335
2336
2337
2338
2339
2340
2341
2342
2343
2344
2345
2346
2347
2348
2349
2350
2351
2352
2353
2354
2355
2356
2357
2358
2359
2360
2361
2362
2363
2364
2365
2366
2367
2368
2369
2370
2371
2372
2373
2374
2375
2376
2377
2378
2379
2380
2381
2382
2383
2384
2385
2386
2387
2388
2389
2390
2391
2392
2393
2394
2395
2396
2397
2398
2399
2400
2401
2402
2403
2404
2405
2406
2407
2408
2409
2410
2411
2412
2413
2414
2415
2416
2417
2418
2419
2420
2421
2422
2423
2424
2425
2426
2427
2428
2429
2430
2431
2432
2433
2434
2435
2436
2437
2438
2439
2440
2441
2442
2443
2444
2445
2446
2447
2448
2449
2450
2451
2452
2453
2454
2455
2456
2457
2458
2459
2460
2461
2462
2463
2464
2465
2466
2467
2468
2469
2470
2471
2472
2473
2474
2475
2476
2477
2478
2479
2480
2481
2482
2483
2484
2485
2486
2487
2488
2489
2490
2491
2492
2493
2494
2495
2496
2497
2498
2499
2500
2501
2502
2503
2504
2505
2506
2507
2508
2509
2510
2511
2512
2513
2514
2515
2516
2517
2518
2519
2520
2521
2522
2523
2524
2525
2526
2527
2528
2529
2530
2531
2532
2533
2534
2535
2536
2537
2538
2539
2540
2541
2542
2543
2544
2545
2546
2547
2548
2549
2550
2551
2552
2553
2554
2555
2556
2557
2558
2559
2560
2561
2562
2563
2564
2565
2566
2567
2568
2569
2570
2571
2572
2573
2574
2575
2576
2577
2578
2579
2580
2581
2582
2583
2584
2585
2586
2587
2588
2589
2590
2591
2592
25
```

## Vamos por el primero... Modelo de regresion

```
In [121]: 1 # Fitting the Logistic regression model
          2 lg = LogisticRegression()
          3
          4 lg.fit(x_train,y_train)
```

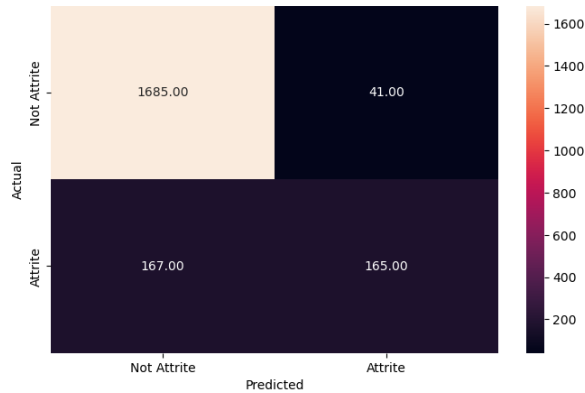
Out[121]: LogisticRegression()

In a Jupyter environment, please rerun this cell to show the HTML representation or trust the notebook.  
On GitHub, the HTML representation is unable to render, please try loading this page with nbviewer.org.

### Revisemos el rendimiento de nuestro modelo

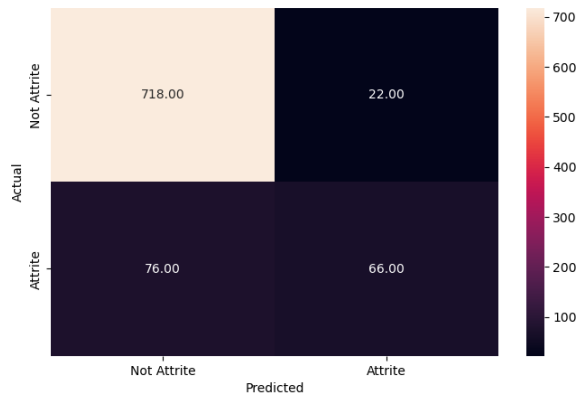
```
In [122]: 1 # Evaluando el rendimiento en nuestro training data
          2 y_pred_train = lg.predict(x_train)
          3 metrics_score(y_train, y_pred_train)
```

	precision	recall	f1-score	support
0	0.91	0.98	0.94	1726
1	0.80	0.50	0.61	332
accuracy			0.90	2058
macro avg	0.86	0.74	0.78	2058
weighted avg	0.89	0.90	0.89	2058



```
In [123]: 1 # Evaluando el rendimiento en nuestro set data
          2 y_pred_test = lg.predict(x_test)
          3 metrics_score(y_test, y_pred_test)
```

	precision	recall	f1-score	support
0	0.90	0.97	0.94	740
1	0.75	0.46	0.57	142
accuracy			0.89	882
macro avg	0.83	0.72	0.76	882
weighted avg	0.88	0.89	0.88	882



### Observaciones:

1. La primera observación es que el modelo tiene una precisión de alrededor del 90% tanto en el conjunto de datos de entrenamiento como en el de prueba. Esto significa que el modelo predice correctamente si un empleado abandonará o no la empresa en el 90% de los casos.
2. La segunda observación es que la recuperación de la clase 1 (empleados que abandonan la empresa) sólo ronda el 50% en los datos de entrenamiento y el 46% en los datos de prueba. La recuperación mide la eficacia del modelo a la hora de identificar a los empleados que corren el riesgo de abandonar la empresa. Una recuperación baja significa que el modelo no es bueno para identificar a los empleados que corren el riesgo de abandonar la empresa.
3. La tercera observación es que, como la recuperación es baja, el modelo no funcionará bien a la hora de distinguir a los empleados que tienen más probabilidades de abandonar la empresa. Esto significa que no ayudará a reducir la tasa de abandono.
4. La cuarta observación es que, como podemos ver en la matriz de confusión, este modelo no es bueno para identificar a los empleados que corren el riesgo de abandonar la empresa. Una matriz de confusión es una tabla que muestra el rendimiento de un modelo comparando sus predicciones con los resultados reales.

En resumen, aunque el modelo tiene una alta precisión, su baja recuperación significa que no es bueno para identificar a los empleados que corren el riesgo de abandonar la empresa. Esto significa que puede no ser eficaz para reducir las bajas.

Comprobemos los coeficientes y averigüemos qué variables provocan el abandono y cuáles pueden ayudar a reducirlo.

```
In [124]: 1 # Visualizando el coeficiente de la regresion logistica
2 cols = X.columns
3
4 coef_lg = lg.coef_
5
6 pd.DataFrame(coef_lg.columns = cols).T.sort_values(by = 0, ascending = False)
```

Out[124]:

	0
OverTime	0.954859
BusinessTravel_Travel_Frequently	0.720271
MaritalStatus_Single	0.623614
YearsSinceLastPromotion	0.551366
YearsAtCompany	0.534303
NumCompaniesWorked	0.496154
BusinessTravel_Travel_Rarely	0.446879
DistanceFromHome	0.386476
JobRole_Human Resources	0.386061
EducationField_Technical Degree	0.300268
MaritalStatus_Married	0.291952
JobLevel_5	0.272820
EducationField_Marketing	0.219969
Gender_Male	0.166946
Education_3	0.159227
EducationField_Human Resources	0.135193
Education_2	0.133188
Education_4	0.115034
JobRole_Marketing Analyst	0.096392
Education_5	0.091082
MonthlyRate	0.058916
EducationField_Programming	0.052074
HourlyRate	0.047247
EducationField_Other	0.035985
JobLevel_3	-0.001782
PerformanceRating	-0.034693
Gamers_Yes	-0.035068
Department_Technology Department	-0.048424
JobRole_Manufacturing Director	-0.072903
PercentSalaryHike	-0.074232
JobRole_Developer	-0.074448
DailyRate	-0.093721
StockOptionLevel	-0.103986
JobLevel_4	-0.169284
JobRole_Manager	-0.183561
JobRole_Customer Experience	-0.186324
WorkLifeBalance	-0.210766
TrainingTimesLastYear	-0.242308
Age	-0.271521
RelationshipSatisfaction	-0.313403
JobSatisfaction	-0.372680
YearsWithCurrManager	-0.391109
YearsInCurrentRole	-0.442057
EnvironmentSatisfaction_2	-0.445452
JobInvolvement_2	-0.480720
TotalWorkingYears	-0.490906
Department_Human Resources	-0.495187
JobRole_Designer	-0.496220
JobRole_Research Director	-0.498261
EnvironmentSatisfaction_3	-0.502891
MonthlyIncome	-0.615774
JobInvolvement_4	-0.637248
EnvironmentSatisfaction_4	-0.652362
JobLevel_2	-0.727879
JobInvolvement_3	-0.744691

Aqui podemos observar que los que estan mas arriba y son de numeros positivos corresponden a los que afectan positivamente al % de salidas de colaboradores en RiotGames.

Por el contrario, los que se encuentran abajo, con numeros negativos, son los que afectan negativamente al nivel de salidas de colaboradores.

Los coeficientes del modelo de regresión logística nos dan el logaritmo de probabilidades, que es difícil de interpretar en el mundo real. Podemos convertir el logaritmo de probabilidades en probabilidades tomando su exponencial.

```
In [125]: 1 odds = np.exp(lg.coef_[0]) # Transformando y encontrando Las probabilidades
2
3 #Agregando Las probabilidades de nuestra base de datos y ordenandolas
4 pd.DataFrame(odds, x_train.columns, columns = ['odds']).sort_values(by = 'odds', ascending = False)
```

Out[125]:

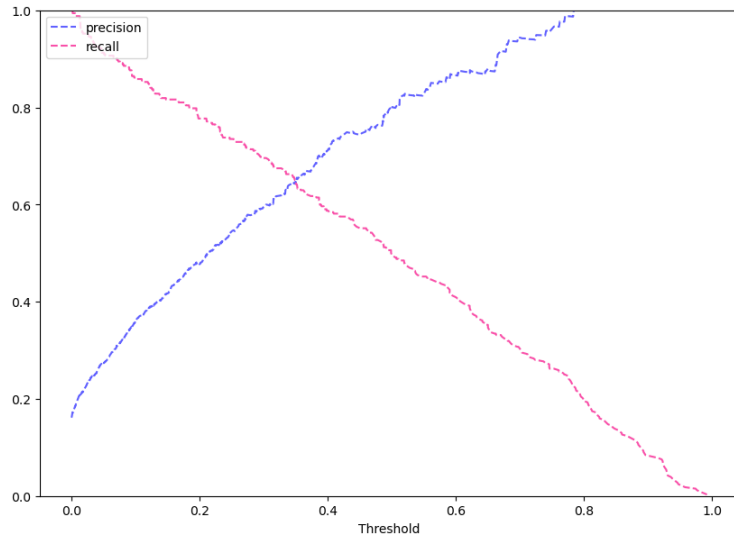
	odds
OverTime	2.598305
BusinessTravel_Travel_Frequently	2.054990
MaritalStatus_Single	1.865658
YearsSinceLastPromotion	1.735622
YearsAtCompany	1.706259
NumCompaniesWorked	1.642392
BusinessTravel_Travel_Rarely	1.563426
DistanceFromHome	1.471786
JobRole_Human Resources	1.471175
EducationField_Technical Degree	1.350220
MaritalStatus_Married	1.339039
JobLevel_5	1.313663
EducationField_Marketing	1.246038
Gender_Male	1.181691
Education_3	1.172605
EducationField_Human Resources	1.144757
Education_2	1.142464
Education_4	1.121911
JobRole_Marketing Analyst	1.101191
Education_5	1.095359
MonthlyRate	1.060886
EducationField_Programming	1.053453
HourlyRate	1.048381
EducationField_Other	1.036640
JobLevel_3	0.998220
PerformanceRating	0.965902
Gamers_Yes	0.965540
Department_Technology Department	0.952730
JobRole_Manufacturing Director	0.929691
PercentSalaryHike	0.928456
JobRole_Developer	0.928256
DailyRate	0.910536
StockOptionLevel	0.901238
JobLevel_4	0.844269
JobRole_Manager	0.832301
JobRole_Customer Experience	0.830004
WorkLifeBalance	0.809963
TrainingTimesLastYear	0.784814
Age	0.762219
RelationshipSatisfaction	0.730955
JobSatisfaction	0.688886
YearsWithCurrManager	0.676306
YearsInCurrentRole	0.642713
EnvironmentSatisfaction_2	0.640535
JobInvolvement_2	0.618338
TotalWorkingYears	0.612071
Department_Human Resources	0.609457
JobRole_Designer	0.608828
JobRole_Research Director	0.607587
EnvironmentSatisfaction_3	0.604780
MonthlyIncome	0.540223
JobInvolvement_4	0.528745
EnvironmentSatisfaction_4	0.520814
JobLevel_2	0.482932
JobInvolvement_3	0.474881

Observaciones:

- Las probabilidades de que un empleado que hace horas extras se desgaste son 2,6 veces superiores a las de uno que no hace horas extras, probablemente porque hacer horas extras no es sostenible durante mucho tiempo para ningún empleado, y puede provocar agotamiento e insatisfacción laboral.
- Las probabilidades de que un empleado que viaja con frecuencia se desgaste son el doble que las de un empleado que no viaja tan a menudo.
- Las probabilidades de que los empleados solteros abandonen el trabajo son aproximadamente 1,85 veces (un 85% más altas) que las de un empleado con otro estado civil.

Curva de precisión/recuperación de la regresión logística

```
In [126]: 1 y_scores_lg = lg.predict_proba(x_train) # predict_proba nos da la probabilidad de cada observacion belonging to each class
2
3
4 precisions_lg, recalls_lg, thresholds_lg = precision_recall_curve(y_train, y_scores_lg[:, 1])
5
6 # Plot values of precisions, recalls, and thresholds
7 plt.figure(figsize = (10, 7))
8
9 plt.plot(thresholds_lg, precisions_lg[:-1], 'b--', label = 'precision', color = '#6161ff')
10
11 plt.plot(thresholds_lg, recalls_lg[:-1], 'g--', label = 'recall', color = '#f950a9')
12
13 plt.xlabel('Threshold')
14
15 plt.legend(loc = 'upper left')
16
17 plt.ylim([0, 1])
18
19 plt.show()
```



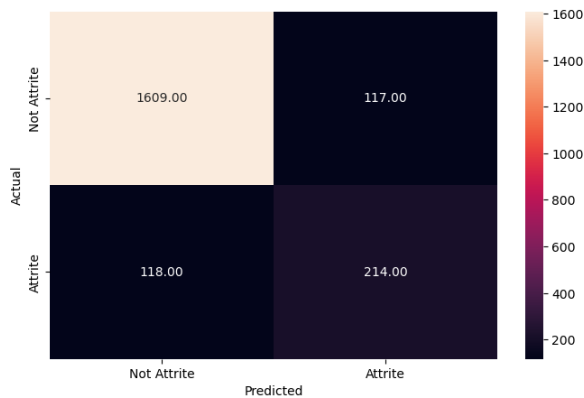
#### Observación:

- Podemos ver que la precisión y el recuerdo están equilibrados para un umbral de alrededor de 0,35.
- El threshold determina el punto en el que el modelo clasifica una observación como positiva o negativa. Al ajustar este valor, puede controlar el equilibrio entre precisión y recuperación. Un umbral más bajo dará lugar a un mayor número de predicciones positivas, lo que aumentará la recuperación pero reducirá potencialmente la precisión. Un umbral más alto dará lugar a menos predicciones positivas, lo que aumentará la precisión pero podría reducir la recuperación.
- En este caso, la observación indica que cuando el umbral se fija en torno a 0,35, los valores de precisión y recuperación (Recall) se equilibran. Esto significa que, con este valor de umbral, el modelo es capaz de identificar correctamente una buena proporción de los casos positivos reales, al tiempo que mantiene un alto nivel de precisión en sus predicciones positivas.

Averiguemos el rendimiento del modelo en este umbral.

```
In [127]: 1 optimal_threshold1 = .35
2
3 y_pred_train = lg.predict_proba(x_train)
4
5 metrics_score(y_train, y_pred_train[:, 1] > optimal_threshold1)
```

	precision	recall	f1-score	support
0	0.93	0.93	0.93	1726
1	0.65	0.64	0.65	332
accuracy			0.89	2058
macro avg	0.79	0.79	0.79	2058
weighted avg	0.89	0.89	0.89	2058



#### Observaciones:

- El rendimiento del modelo ha mejorado. La recuperación ha aumentado significativamente para la clase 1.

Comprobemos el rendimiento con los datos de prueba (set data).

```
In [128]: 1 optimal_threshold1 = .35
          2
          3 y_pred_test = lg.predict_proba(x_test)
          4
          5 metrics_score(y_test, y_pred_test[:, 1] > optimal_threshold1)
```

	precision	recall	f1-score	support
0	0.93	0.92	0.92	740
1	0.60	0.62	0.61	142
accuracy			0.87	882
macro avg	0.76	0.77	0.77	882
weighted avg	0.87	0.87	0.87	882



**Observaciones:**

- El modelo ofrece un rendimiento similar en los conjuntos de datos de prueba y de entrenamiento, es decir, ofrece un rendimiento generalizado.
- La recuperación de los datos de prueba ha aumentado y, al mismo tiempo, la precisión ha disminuido ligeramente, como era de esperar al ajustar el umbral.
- La recuperación y la precisión medias del modelo son buenas, pero veamos si podemos obtener un rendimiento aún mejor utilizando otros algoritmos.

**K-Nearest Neighbors (K-NN)**

K-NN utiliza características de los datos de entrenamiento para predecir los valores de nuevos puntos de datos, lo que significa que al nuevo punto de datos se le asignará un valor basado en lo similar que es a los puntos de datos del conjunto de entrenamiento.

¿Que debemos hacer?

1. Seleccionar K
2. Calcular la distancia (Euclídea, Manhattan, etc.)
3. Encontrar los K vecinos más cercanos
4. Votar por mayoría las etiquetas

La "K" en el algoritmo K-NN es el número de vecinos más cercanos entre los que queremos votar. Generalmente, K es un número impar cuando el número de clases es par, para obtener un voto mayoritario. Supongamos que K=3. En ese caso, haremos un círculo con el nuevo punto de datos como centro tan grande como encerrar sólo los tres puntos de datos más cercanos en el plano.

**Pero antes de construir realmente el modelo, necesitamos identificar el valor de K que se utilizará en K-NN. Para ello realizaremos los siguientes pasos.**

- Para cada valor de K (de 1 a 15), dividir el conjunto de entrenamiento en un nuevo conjunto de entrenamiento y validación (30 veces).
- Escalar los datos de entrenamiento y los datos de validación
- Calcula la media del error en los conjuntos de entrenamiento y validación para cada valor de K.
- Represente gráficamente el error medio de entrenamiento frente al de validación para todos los K.
- Elegir el valor óptimo de K en el gráfico para que los dos errores sean comparables.



```

In [139]: 1 knn = KNeighborsClassifier()
2
3 # We select the optimal value of K for which the error rate is the Least in the validation data
4 # Let us Loop over a few values of K to determine the optimal value of K
5
6 train_error = []
7
8 test_error = []
9
10 knn_many_split = {}
11
12 error_df_knn = pd.DataFrame()
13
14 features = X.columns
15
16 for k in range(1, 15):
17     train_error = []
18     test_error = []
19
20     lista = []
21
22     knn = KNeighborsClassifier(n_neighbors = k)
23
24     for i in range(30):
25         x_train_new, x_val, y_train_new, y_val = train_test_split(x_train, y_train, test_size = 0.20)
26
27         # Fitting K-NN on the training data
28         knn.fit(x_train_new, y_train_new)
29
30         # Calculating error on the training data and the validation data
31         train_error.append(1 - knn.score(x_train_new, y_train_new))
32
33         test_error.append(1 - knn.score(x_val, y_val))
34
35     lista.append(sum(train_error)/len(train_error))
36
37     lista.append(sum(test_error)/len(test_error))
38
39     knn_many_split[k] = lista
40
41 knn_many_split
42

```

```

Out[139]: {1: [0.0, 0.09555016181229771],
2: [0.05516403402187121, 0.14579288025889972],
3: [0.05838396111786149, 0.15800970873786416],
4: [0.11814499797488862, 0.15210355987055021],
5: [0.109437019036047, 0.1519417475728156],
6: [0.13063993519643577, 0.15056634304207123],
7: [0.1228230052652896, 0.15307443365695797],
8: [0.13681652490887, 0.15420711974110032],
9: [0.13140947752126367, 0.1545307443365696],
10: [0.1445727014985824, 0.15024271844660203],
11: [0.13997569866342646, 0.15064724919093855],
12: [0.14746861077359255, 0.15574433656957934],
13: [0.14400567031186715, 0.14595469255663435],
14: [0.14868367760226814, 0.15485436893203888]}

```

**PROYECTO EN PROCESO. TODAVÍA NO FINALIZADO.**