

# Informe Trabajo Integrador

## Análisis de datos

CEIA 2021

## Contenido

Introducción: .....	3
Análisis exploratorio inicial .....	3
Visualización de las primeras filas del dataset .....	3
Detalle de variables .....	3
Identificación de tipos de datos: .....	4
Análisis por tipo de variable: numérica .....	4
Análisis por tipo de variable: categórica .....	5
Análisis de la variable de salida y balance de clases: .....	5
Preprocesamiento del dataset .....	6
Ejemplo Modelo 1 .....	6
Manejo de datos imbalanceados .....	6
Tratamiento de NaNs .....	7
Ingeniería de features básica .....	7
Ejemplo Modelo 2 .....	8
Tratamiento de outliers .....	8
Tratamiento de NaNs .....	8
Ingeniería de features básica .....	9
Ejemplo Modelo 3 .....	9
Manejo de datos imbalanceados .....	10
Tratamiento de NaNs .....	10
Ingeniería de features básica .....	10
Implementación de modelos de ML .....	11
Ejemplo Modelo 1 .....	11
Modelo Base .....	11
Regresión Logística .....	12
Random Forest .....	12
Deep learning .....	13
Ejemplo Modelo 2 .....	13
Modelo Base .....	13
Regresión logística .....	14
XGBoost .....	14
Ejemplo Modelo 3 .....	14
XGBoost .....	14
Referencias .....	15

## Introducción:

A partir de la propuesta de trabajo final para la materia de Análisis de Datos, se elabora el presente informe final con el objetivo de detallar los métodos utilizados para la resolución del trabajo integrador.

El dataset brindado para realizar este trabajo proviene de *kaggle* y el objetivo es predecir la variable de salida *RainTomorrow*. El dataset contiene aproximadamente mas de 10 años de observaciones meteorológicas de varios lugares de Australia

El desafío consiste en realizar un análisis exploratorio inicial utilizando diferentes técnicas para el análisis de datos. Luego, se realizará la preparación de los datos que serán el input del modelo de machine learning, el entrenamiento, evaluación e interpretación de resultados de los diferentes modelos. Finalmente se deberán desarrollar las conclusiones.

## Análisis exploratorio inicial

### Visualización de las primeras filas del dataset

	Date	Location	MinTemp	MaxTemp	Rainfall	Evaporation	Sunshine	WindGustDir	WindGustSpeed	WindDir9am	WindDir3pm	WindSpeed9am	WindSpeed3pm	Humidity9am
0	2008-12-01	Albury	13.4	22.9	0.6	NaN	NaN	W	44.0	W	WNW	20.0	24.0	71.0
1	2008-12-02	Albury	7.4	25.1	0.0	NaN	NaN	WNW	44.0	NNW	WSW	4.0	22.0	44.0
2	2008-12-03	Albury	12.9	25.7	0.0	NaN	NaN	WSW	46.0	W	WSW	19.0	26.0	38.0
3	2008-12-04	Albury	9.2	28.0	0.0	NaN	NaN	NE	24.0	SE	E	11.0	9.0	45.0
4	2008-12-05	Albury	17.5	32.3	1.0	NaN	NaN	W	41.0	ENE	NW	7.0	20.0	82.0

### Detalle de variables

- Date: La fecha de la observación
- Location: El nombre común de la ubicación de la estación meteorológica.
- MinTemp: La temperatura mínima en grados centígrados.
- MaxTemp: La temperatura máxima en grados centígrados.
- Rainfall: La cantidad de lluvia registrada para el día en mm
- Evaporation: La denominada evaporación de la bandeja de clase A (mm) en las 24 horas a las 9 am.
- Sunshine: El número de horas de sol brillante en el día.
- WindGustDir: La dirección de la ráfaga de viento más fuerte en las 24 horas hasta la medianoche.
- WindGustSpeed: La velocidad (km / h) de la ráfaga de viento más fuerte en las 24 horas hasta la medianoche.
- WindDir9am: Dirección del viento a las 9 a. M.
- WindDir3pm: Dirección del viento a las 3pm.
- WindSpeed9am: Velocidad del viento (km / h) promediada durante 10 minutos antes de las 9 a. M.
- WindSpeed3pm: Velocidad del viento (km / h) promediada durante 10 minutos antes de las 3 p.m.

- Humidity9am: Humedad (porcentaje) a las 9 a. M.
- Humidity3pm: Humedad (porcentaje) a las 3 p.m.
- Pressure9am: Presión atmosférica (hpa) reducida al nivel medio del mar a las 9 a. M.
- Pressure3pm: Presión atmosférica (hpa) reducida al nivel medio del mar a las 3 p.m.
- Cloud9am: Fracción de cielo oscurecida por nubes a las 9 a. M. Esto se mide en "octas", que son una unidad de octavos.
- Cloud3pm: Fracción de cielo oscurecida por nubes (en "octas": octavos) a las 15.00 horas.
- Temp9am: Temperatura (grados C) a las 9 a. M.
- Temp3pm: Temperatura (grados C) a las 3 p.m.
- RainToday: Booleano: 1 si la precipitación (mm) en las 24 horas hasta las 9 a.m. excede 1 mm, de lo contrario 0
- RainTomorrow: Se utiliza para crear la variable de respuesta RainTomorrow.

El dataset contiene 145460 filas y 24 columnas.

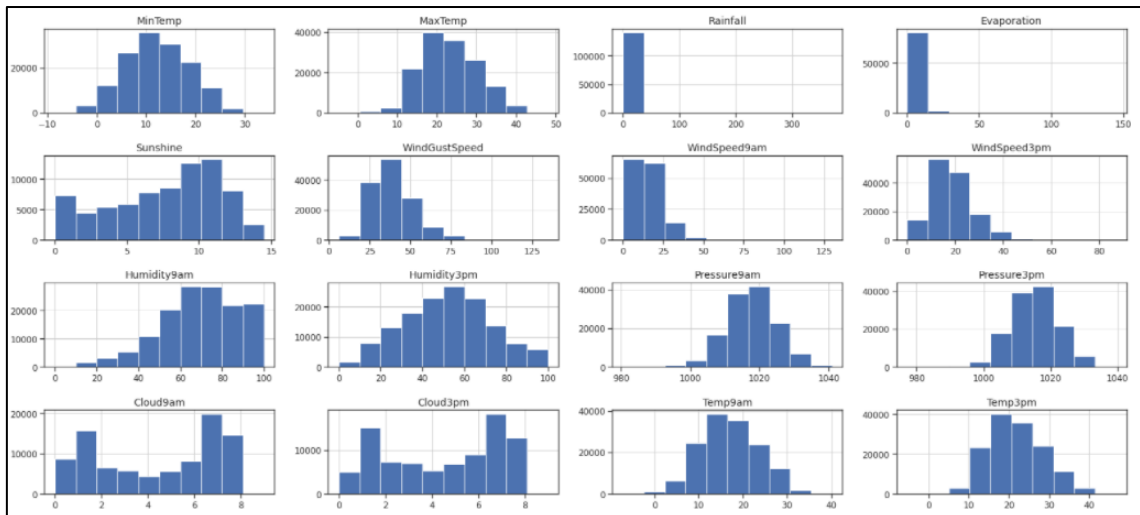
Identificación de tipos de datos:

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 145460 entries, 0 to 145459
Data columns (total 23 columns):
#   Column                Non-Null Count  Dtype
---  ---
0   Date                  145460 non-null object
1   Location              145460 non-null object
2   MinTemp               143975 non-null float64
3   MaxTemp              144199 non-null float64
4   Rainfall              142199 non-null float64
5   Evaporation           82670 non-null float64
6   Sunshine              75625 non-null float64
7   WindGustDir           135134 non-null object
8   WindGustSpeed         135197 non-null float64
9   WindDir9am           134894 non-null object
10  WindDir3pm            141232 non-null object
11  WindSpeed9am          143693 non-null float64
12  WindSpeed3pm          142398 non-null float64
13  Humidity9am           142806 non-null float64
14  Humidity3pm           140953 non-null float64
15  Pressure9am           130395 non-null float64
16  Pressure3pm           130432 non-null float64
17  Cloud9am              89572 non-null float64
18  Cloud3pm              86102 non-null float64
19  Temp9am               143693 non-null float64
20  Temp3pm               141851 non-null float64
21  RainToday             142199 non-null object
22  RainTomorrow          142193 non-null object
dtypes: float64(16), object(7)
memory usage: 25.5+ MB
```

Las variables categóricas en primera instancia son *Date*, *Location*, *WindGustDir*, *WindDir9am*, *WindDir3pm*, *RainToday*. La variable a predecir es *RainTomorrow* (variable de salida) y las demás variables son numéricas. Todas las variables menos la variable a predecir (*RainTomorrow*) son input para el modelo (variables de entrada).

Análisis por tipo de variable: numérica

Se observa la distribución de los datos a partir de la creación de histograma. La mayoría de las variables muestra seguir una distribución normal.



## Análisis por tipo de variable: categórica

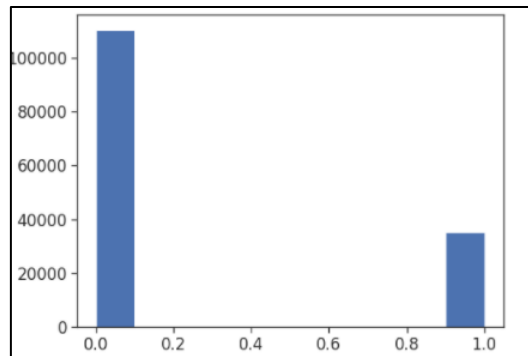
Al realizar el análisis de cardinalidad de las variables categóricas se llega a la conclusión de que Date y Location tienen alta cardinalidad.

Inicialmente durante el análisis se identificó a la variable Date de importancia para el modelo ya que el mes podría estar relacionado con la variable de salida. Más adelante, durante el proceso se desestimó ya que al realizar un análisis de correlación entre el mes extraído de la columna Date y la variable de salida, se pudo concluir que no estaban correlacionadas.

MinTemp	1.00	0.74	0.10	0.47	0.07	0.18	0.18	0.18	-0.23	0.01	-0.45	-0.46	0.08	0.02	0.90	0.71	-0.20	0.08
MaxTemp	0.74	1.00	-0.07	0.59	0.47	0.07	0.01	0.05	-0.50	-0.51	-0.33	-0.43	-0.29	-0.28	0.89	0.98	-0.16	-0.16
Rainfall	0.10	-0.07	1.00	-0.06	-0.23	0.13	0.09	0.06	0.22	0.26	-0.17	-0.13	0.20	0.17	0.01	-0.08	-0.03	0.24
Evaporation	0.47	0.59	-0.06	1.00	0.37	0.20	0.19	0.13	-0.50	-0.39	-0.27	-0.29	-0.18	-0.18	0.55	0.57	-0.03	-0.12
Sunshine	0.07	0.47	-0.23	0.37	1.00	-0.03	0.01	0.05	-0.49	-0.43	0.04	-0.02	-0.68	-0.70	0.29	0.49	0.02	-0.44
WindGustSpeed	0.18	0.07	0.13	0.20	-0.03	1.00	0.81	0.69	-0.22	-0.03	-0.46	-0.41	0.07	0.11	0.15	0.03	0.06	0.23
WindSpeed9am	0.18	0.01	0.09	0.19	0.01	0.61	1.00	0.52	-0.27	-0.03	-0.23	-0.18	0.03	0.05	0.13	0.00	0.05	0.10
WindSpeed3pm	0.18	0.05	0.06	0.13	0.05	0.69	0.52	1.00	-0.15	0.02	-0.30	-0.26	0.05	0.03	0.16	0.03	0.06	0.09
Humidity9am	-0.23	-0.50	0.22	-0.50	-0.49	-0.22	-0.27	-0.15	1.00	0.67	0.14	0.19	0.45	0.36	-0.47	-0.50	-0.09	0.25
Humidity3pm	0.01	-0.51	0.26	-0.39	-0.40	-0.03	-0.03	0.02	0.67	1.00	-0.03	0.05	0.52	0.52	-0.22	-0.56	-0.02	0.43
Pressure9am	-0.45	-0.33	-0.17	-0.21	0.04	-0.46	-0.23	-0.30	0.14	-0.03	1.00	0.75	-0.13	-0.15	-0.12	-0.29	0.03	-0.24
Pressure3pm	-0.46	-0.43	-0.13	-0.26	0.02	-0.41	-0.18	-0.26	0.19	0.05	0.75	1.00	-0.06	-0.08	-0.42	-0.39	0.03	-0.22
Cloud9am	0.08	-0.29	0.20	-0.18	-0.66	0.07	0.03	0.05	0.45	0.52	-0.13	-0.06	1.00	0.60	-0.14	-0.30	-0.01	0.32
Cloud3pm	0.02	-0.28	0.17	-0.18	-0.70	0.11	0.05	0.03	0.36	0.52	-0.15	-0.08	0.60	1.00	-0.13	-0.32	-0.00	0.38
Temp9am	0.90	0.89	0.01	0.55	0.29	0.15	0.13	0.16	-0.47	-0.22	-0.42	-0.47	-0.14	-0.13	1.00	0.86	-0.14	-0.02
Temp3pm	0.71	0.88	-0.08	0.57	0.49	0.03	0.00	0.03	-0.50	-0.56	-0.29	-0.39	-0.30	-0.32	0.86	1.00	-0.18	-0.19
Month	-0.20	-0.16	-0.03	-0.03	0.02	0.06	0.05	0.06	-0.09	-0.02	0.03	0.03	-0.01	-0.00	-0.14	-0.18	1.00	0.01
binary_rain_tomorrow	0.08	-0.16	0.24	-0.12	-0.44	0.23	0.10	0.09	0.25	0.43	-0.24	-0.22	0.32	0.38	-0.02	-0.19	0.01	1.00

## Análisis de la variable de salida y balance de clases:

Se realizó el análisis del balance de clases en la variable de salida aplicando una función lambda donde "No" = 0 y "Si" =1 y se identificó que las clases estaban desbalanceadas, como se muestra en la imagen:



## Preprocesamiento del dataset

### Ejemplo Modelo 1

Se realiza una conversión binaria como técnica de codificación de la variable de salida para poder utilizarla en el entrenamiento y se crea la nueva variable de salida:

`"binary_rain_tomorrow"`

Pressure9am	Pressure3pm	Cloud9am	Cloud3pm	Temp9am	Temp3pm	RainToday	RainTomorrow	binary_rain_tomorrow
1007.7	1007.1	8.0	NaN	16.9	21.8	No	No	0
1010.6	1007.8	NaN	NaN	17.2	24.3	No	No	0
1007.6	1008.7	NaN	2.0	21.0	23.2	No	No	0
1017.6	1012.8	NaN	NaN	18.1	26.5	No	No	0
1010.8	1006.0	7.0	8.0	17.8	29.7	No	No	0

### Manejo de datos imbalanceados

Como se observó en el análisis exploratorio inicial, los datos están fuertemente desbalanceados en favor de los casos de que no llueva al día siguiente. Frente a esto se pueden adoptar distintas estrategias:

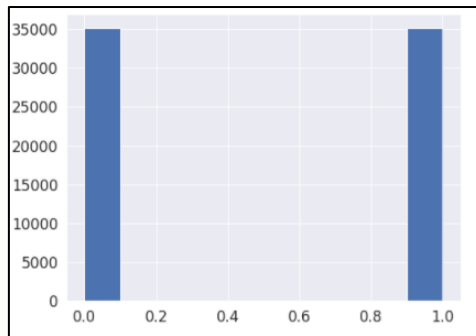
Upsampling: aumentar la cantidad de muestras de la clase minoritaria (agregando muestras reales o artificiales).

Downsampling: reducir la cantidad de muestras de la clase mayoritaria.

Seleccionar cuidadosamente la métrica de evaluación: no hacer nada en la preparación del dataset, y elegir una métrica de performance adecuada para este escenario (TPR/TNR/AUC/Precision/Recall/etc.)

Para continuar con este análisis trabajando con datos originales se procederá con la segunda opción.

Se genera un nuevo dataset que contiene clases balanceadas y se verifica que realmente quedaron balanceadas:



El nuevo dataset “df\_downsampled” contiene 70288 filas y 24 columnas.

### Tratamiento de NaNs

Se identifica un gran porcentaje de NaNs:

```
Cantidad de filas con nans (89040, 23)
Cantidad de filas sin nans (56420, 23)
```

Se realiza un análisis de la cantidad de NaNs del nuevo dataset y el tratamiento de los mismos.

Para el tratamiento de NaNs se utiliza la **imputación estadística** con la cual se completaron las variables categóricas con la moda y las variables numéricas con la media.

Se verifica que las distribuciones no hayan cambiado respecto al dataset original y que ya no haya filas con NaNs:

```
Cantidad de filas con nans (0, 24)
Cantidad de filas sin nans (70288, 24)
```

### Ingeniería de features básica

- Para la variable *RainToday* se identifica que solo tiene dos clases, por lo cual se realiza una transformación con one hot encoding, generando dos nuevas columnas: *raintoday\_No* y *raintoday\_Yes*.

#	Column	Non-Null Count	Dtype
0	Date	70288 non-null	object
1	Location	70288 non-null	object
2	MinTemp	70288 non-null	float64
3	MaxTemp	70288 non-null	float64
4	Rainfall	70288 non-null	float64
5	Evaporation	70288 non-null	float64
6	Sunshine	70288 non-null	float64
7	WindGustDir	70288 non-null	object
8	WindGustSpeed	70288 non-null	float64
9	WindDir9am	70288 non-null	object
10	WindDir3pm	70288 non-null	object
11	WindSpeed9am	70288 non-null	float64
12	WindSpeed3pm	70288 non-null	float64
13	Humidity9am	70288 non-null	float64
14	Humidity3pm	70288 non-null	float64
15	Pressure9am	70288 non-null	float64
16	Pressure3pm	70288 non-null	float64
17	Cloud9am	70288 non-null	float64
18	Cloud3pm	70288 non-null	float64
19	Temp9am	70288 non-null	float64
20	Temp3pm	70288 non-null	float64
21	RainToday	70288 non-null	object
22	RainTomorrow	70288 non-null	object
23	Date_dat	70288 non-null	datetime64[ns]
24	Month	70288 non-null	int64
25	binary_rain_tomorrow	70288 non-null	int64
26	raintoday_No	70288 non-null	uint8
27	raintoday_Yes	70288 non-null	uint8

dtypes: datetime64[ns](1), float64(16), int64(2), object(7), uint8(2)  
memory usage: 14.6+ MB  
None

- Para las demás variables categóricas se hizo el mismo tratamiento, pero al analizar la correlación con la variable de salida, éstas no eran significativas, por lo cual se decidió omitirlas y excluirlas del dataset.

## Ejemplo Modelo 2

### Tratamiento de outliers

Se crearon funciones para encontrar los outliers dentro del dataset y se identificaron los límites superior e inferior de las variables: Rainfall, Evaporation, WindGustSpeed, WindSpeed9am y WindSpeed3pm. Se eliminaron los valores atípicos para que los mismos no afecten a los modelos de clasificación.

### Tratamiento de NaNs

Se identificó un gran porcentaje de NaNs, y para tratarlos se utilizó el método de imputación estadística, completando los nans de las variables categóricas con la moda y los de las variables numéricas con la mediana.

Grafico antes de aplicar el método de imputación estadística:



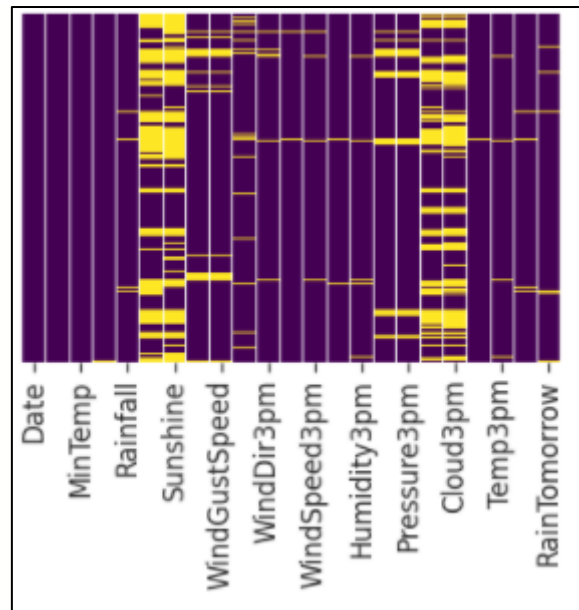
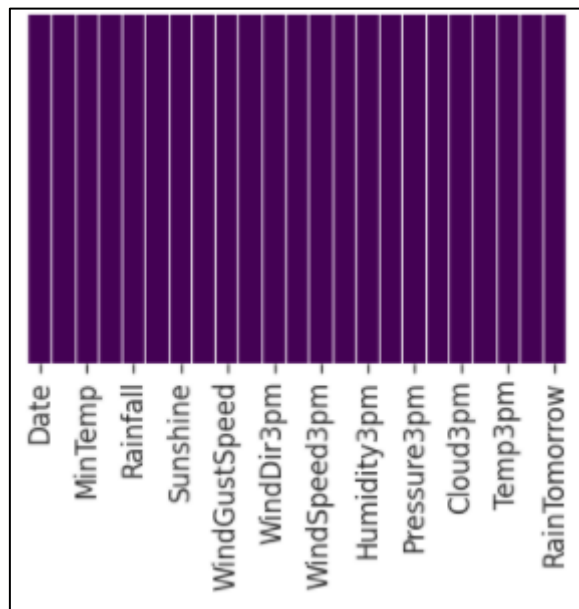


Grafico luego de aplicar el método de imputación estadística (mediana y moda):



### Ingeniería de features básica

Se utiliza el método de one hot encoding para todas las variables categóricas, menos para la variable Date, la cual se excluye del análisis.

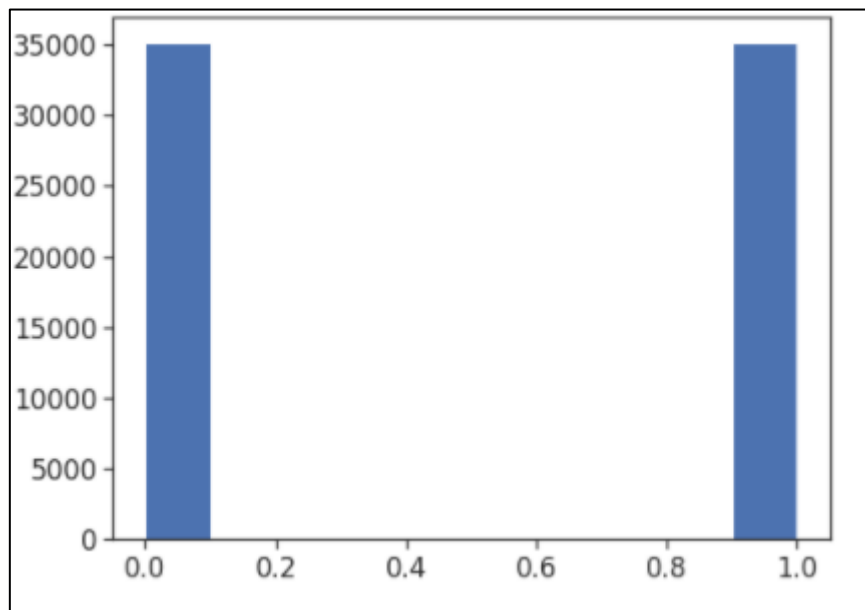
Para las variables numéricas, se realiza la normalización para la aplicación de los diferentes modelos de clasificación.

### Ejemplo Modelo 3

Para el preprocesamiento de datos, se crea la variable `binary_rain_tomorrow` utilizando la función lambda.

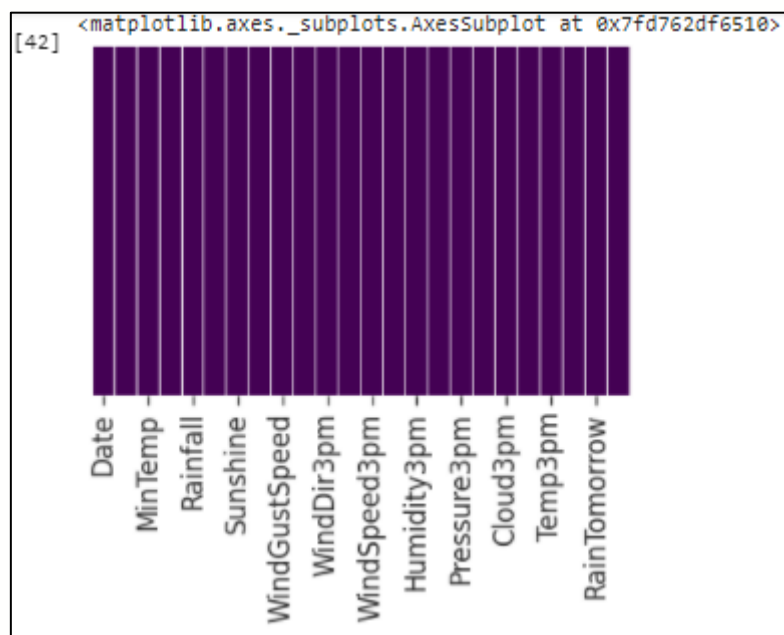
### Manejo de datos imbalanceados

Se aplica la técnica de downsampling y se genera un nuevo dataset "df\_downsampled" para el balance de las clases.



### Tratamiento de NaNs

Se completan los datos categóricos con la moda y los datos numéricos con la media, obteniendo un dataset final sin nans:



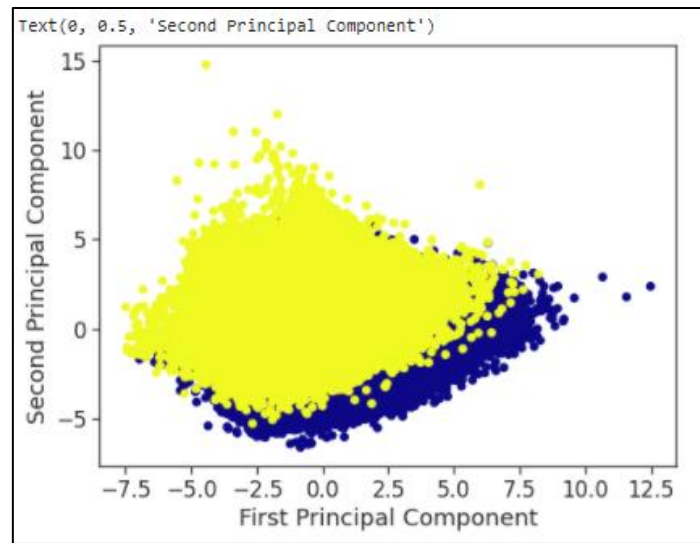
### Ingeniería de features básica

Se elimina la variable Date para seguir con el desarrollo, ya que en el primer modelo implementado se definió que no aporta valor al análisis.

Se realiza la transformación de las variables categoricas utilizando one hot encoding, lo cual genera una gran cantidad de columnas adicionales.

Se realiza la normalización de las variables numericas con `StandardScaler()`

Se utiliza PCA como método de reducción de dimensionalidad:



## Implementación de modelos de ML

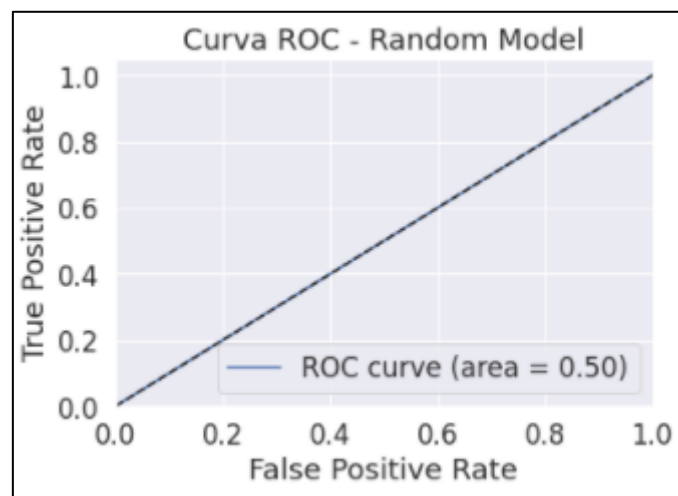
Como métrica de evaluación de los modelos se utilizó AUC porque mide la capacidad del modelo para predecir una mayor puntuación para ejemplos positivos en comparación con ejemplos negativos.

### Ejemplo Modelo 1

#### Modelo Base

Se crea un modelo base a partir de una clase con fit y transform que devuelve resultados completamente random. El mismo se utilizará como instrumento de comparación de los siguientes modelos desarrollados.

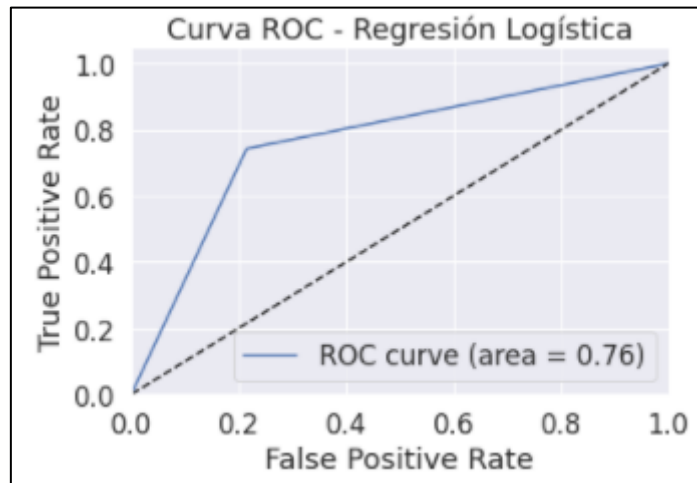
AUC: 0.5



## Regresión Logística

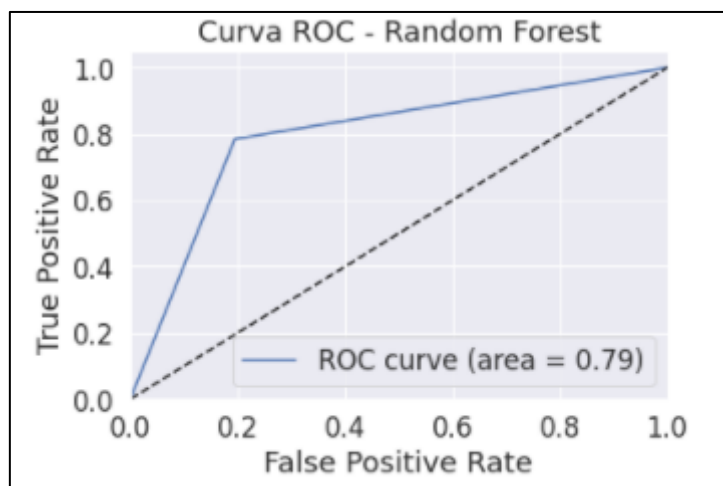
El segundo modelo implementado es la regresión logística, ya que este método trata de modelar la probabilidad de una variable cualitativa binaria (dos posibles valores), en este caso “llueve mañana” o “no llueve mañana”, en función de una o mas variables independientes. Este método se debe aplicar con datos normalizados, quiere decir que los valores de las variables de entrada tienen que oscilar entre 0 y 1. Es por eso que en un paso anterior se utiliza `MinMaxScaler()`

AUC: 0.76



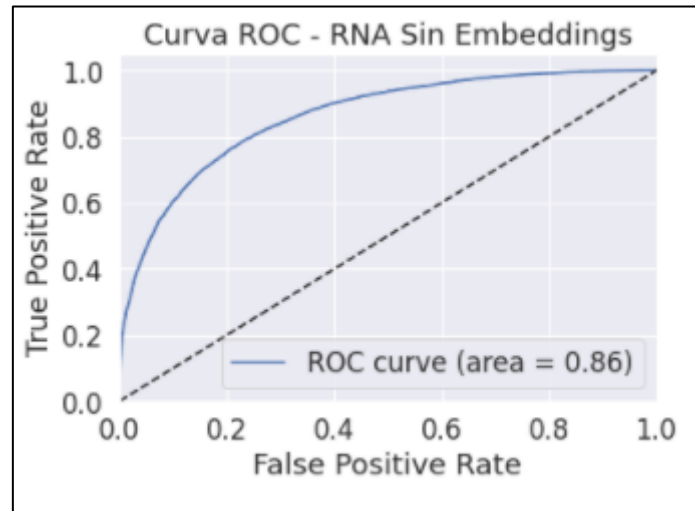
## Random Forest

El tercer modelo implementado es el de *random forest*, un algoritmo de clasificación conformado por muchos árboles de decisión. Utiliza el método de ensamble y la aleatoriedad de las características al construir cada árbol para crear el *random forest* de árboles no correlacionados cuya predicción conjunta sea más precisa que la de cualquier árbol individual. Se obtuvo un AUC de 0.79



## Deeplearning

El ultimo modelo utilizado es el de *deeplearning*. Se utilizo un solo *hidden layer* y datos normalizados. Como función de activación en la salida se utilizó *sigmoid* ya que, al realizar un análisis de la variable de salida, se sabia que la misma iba a ser 0 o 1.



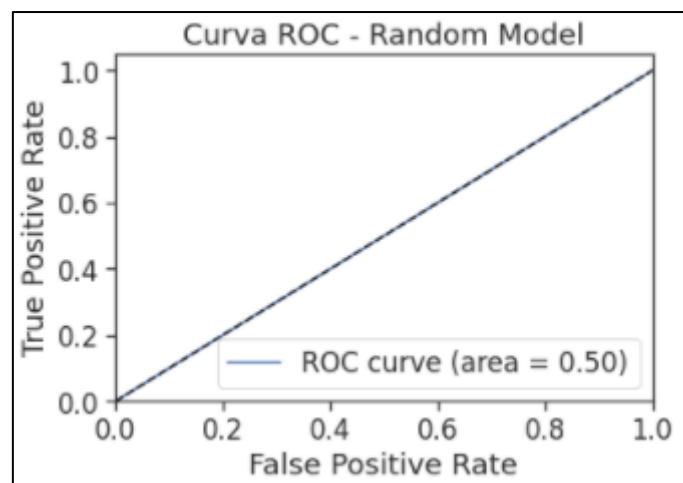
El modelo de DL (con un AUC del 86.16%) tiene un mejor área bajo la curva, no presenta baches por lo cual es un indicador de que clasifica bien (es una curva suave).

## Ejemplo Modelo 2

### Modelo Base

Se crea un modelo base a partir de una clase con fit y transform que devuelve resultados completamente random. El mismo se utilizará como instrumento de comparación de los siguientes modelos desarrollados.

AUC: 0.5



## Regresión logística

El segundo modelo implementado es la regresión logística. Este método se debe aplicar con datos normalizados, quiere decir que los valores de las variables de entrada tienen que oscilar entre 0 y 1. Es por eso que en un paso anterior se utiliza `MinMaxScaler()`

Accuracy: 85 %

### Regresion logistica

```
[67] from sklearn.metrics import confusion_matrix
      from sklearn.metrics import accuracy_score
      model = LogisticRegression(max_iter=500)
      model.fit(X_train, y_train)
      predicted=model.predict(X_test)

      conf = confusion_matrix(y_test, predicted)
      print ("The accuracy of Logistic Regression is : ", accuracy_score(y_test, predicted)*100, "%")

The accuracy of Logistic Regression is : 84.91681561941428 %
```

## XGBoost

El tercer modelo implementado es XGBoost, el cual arrojó mejores resultados que la regresión logística:

```
Accuracy = 0.8614510289197489
ROC score = 0.8930273284320036
```

	precision	recall	f1-score	support
0	0.88	0.95	0.91	33987
1	0.76	0.55	0.64	9651
accuracy			0.86	43638
macro avg	0.82	0.75	0.78	43638
weighted avg	0.85	0.86	0.85	43638

## Ejemplo Modelo 3

### XGBoost

Se implemento el modelo XGBoost el cual arrojo un accuracy de 79%:

```
import xgboost as xgb
xgb = xgb.XGBClassifier()
xgb.fit(X_train, y_train)
y_pred = xgb.predict(X_test)

from sklearn.metrics import mean_squared_error
rmse = np.sqrt(mean_squared_error(y_test, y_pred))
print("RMSE: %f" % (rmse))

RMSE: 0.464908

from sklearn.metrics import accuracy_score
accuracy = accuracy_score(y_test, y_pred)
print("Accuracy: %.2f" % (accuracy * 100.0))

Accuracy: 78.39
```

## Referencias

Ejemplo Modelo 1: ADD- Trabajo Integrador- Modelo 1.ipynb

Ejemplo Modelo 2: ADD- Trabajo Integrador- Modelo 2.ipynb

Ejemplo Modelo 3: ADD- Trabajo Integrador- Modelo 3.ipynb