

TP Integrador – Análisis de Datos

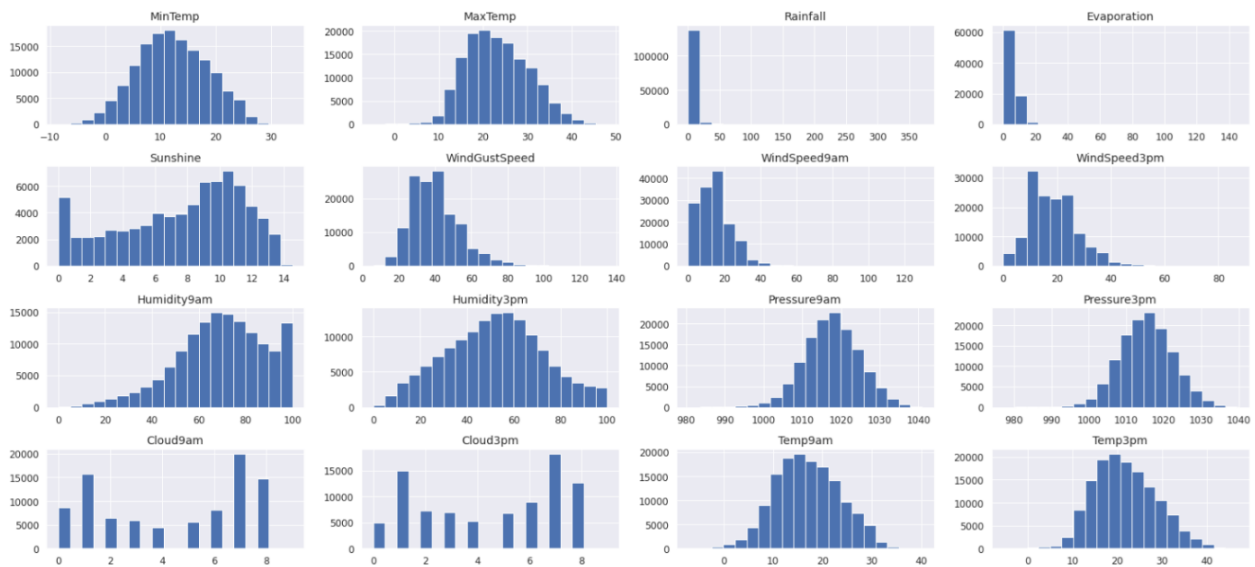
Alumno: Tinelli, Francisco

[Link Google Colab](#) – [Link Dataset Kaggle](#) – [Link del Repositorio Github](#)

1. Análisis exploratorio inicial

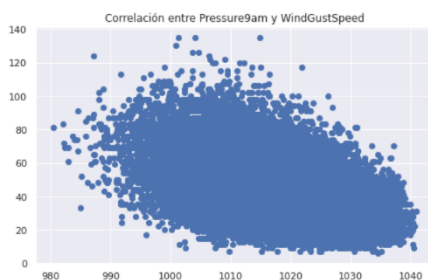
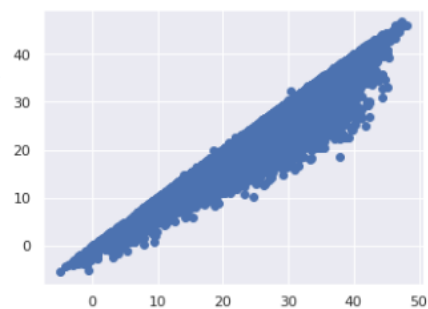
En esta primera etapa, se analizó el dataset, la cantidad de variables presentes y sus tipos (numéricas, categóricas y compuestas) y la variable de salida.

Variables numéricas: En el caso de las variables numéricas, se analizaron mediante histogramas sus distribuciones, algunas de las cuales se asemejan a una normal o a una exponencial.



Se analizó también la correlación entre algunas variables mediante la matriz de correlación y gráficas de correlación.

MinTemp	1	0.74	0.1	0.47	0.073	0.18	0.18	0.18	-0.23	0.0061	-0.45	-0.46	0.079	0.022	0.9	0.71
MaxTemp	0.74	1	-0.075	0.59	0.47	0.068	0.014	0.05	-0.5	-0.51	-0.33	-0.43	-0.29	-0.28	0.89	0.98
Rainfall	0.1	-0.075	1	-0.064	-0.23	0.13	0.087	0.058	0.22	0.26	-0.17	-0.13	0.2	0.17	0.011	-0.08
Evaporation	0.47	0.59	-0.064	1	0.37	0.2	0.19	0.13	-0.5	-0.39	-0.27	-0.29	-0.18	-0.18	0.55	0.57
Sunshine	0.073	0.47	-0.23	0.37	1	-0.035	0.0055	0.054	-0.49	-0.63	0.042	-0.02	-0.68	-0.7	0.29	0.49
WindGustSpeed	0.18	0.068	0.13	0.2	-0.035	1	0.61	0.69	-0.22	-0.026	-0.86	-0.41	0.072	0.11	0.15	0.033
WindSpeed9am	0.18	0.014	0.087	0.19	0.0055	0.61	1	0.52	-0.27	-0.032	-0.23	-0.18	0.025	0.055	0.13	0.0046
WindSpeed3pm	0.18	0.05	0.058	0.13	0.054	0.69	0.52	1	-0.15	0.016	-0.3	-0.26	0.053	0.05	0.16	0.028
Humidity9am	-0.23	-0.5	0.22	-0.5	-0.49	-0.22	-0.27	-0.15	1	0.67	0.14	0.19	0.45	0.36	-0.47	-0.5
Humidity3pm	0.0061	-0.51	0.26	-0.39	-0.63	-0.026	-0.032	0.016	0.67	1	-0.028	0.052	0.52	0.52	-0.22	-0.56
Pressure9am	-0.45	-0.33	-0.17	-0.27	0.042	-0.46	-0.23	-0.3	0.14	-0.028	1	0.96	-0.13	-0.15	-0.42	-0.29
Pressure3pm	-0.46	-0.43	-0.13	-0.29	-0.02	-0.41	-0.18	-0.26	0.19	0.052	0.96	1	-0.061	-0.085	-0.47	-0.39
Cloud9am	0.079	-0.29	0.2	-0.18	-0.68	0.072	0.025	0.053	0.45	0.52	-0.13	-0.061	1	0.6	-0.14	-0.3
Cloud3pm	0.022	-0.28	0.17	-0.18	-0.7	0.11	0.055	0.025	0.36	0.52	-0.15	-0.085	0.6	1	-0.13	-0.32
Temp9am	0.9	0.89	0.011	0.55	0.29	0.15	0.13	0.16	-0.47	-0.22	-0.42	-0.47	-0.14	-0.13	1	0.86
Temp3pm	0.71	0.98	-0.08	0.57	0.49	0.033	0.0046	0.028	-0.5	-0.56	-0.29	-0.39	-0.3	-0.32	0.86	1
MinTemp																
MaxTemp																
Rainfall																
Evaporation																
Sunshine																
WindGustSpeed																
WindSpeed9am																
WindSpeed3pm																
Humidity9am																
Humidity3pm																
Pressure9am																
Pressure3pm																
Cloud9am																
Cloud3pm																
Temp9am																
Temp3pm																



Variables categóricas: Se analizaron los diferentes valores que pueden tomar las variables categóricas: nombres de ciudades, direcciones de viento (N, NE, etc.) y variables True/False.

Variables compuestas: Es el caso de la fecha para la cual se estudió el rango.

Variable de salida: Es la última columna del dataset '*RainTomorrow*' la cual es booleana y está desbalanceada presentando aproximadamente un 76% de casos *False*, 22% *True* y 2% *NaN*.

2. Esquema de validación de resultados

Se eliminaron las muestras con salida NaN y se particionó el dataset en un 70% para entrenamiento y 30% para test.

3. Limpieza y preparación de datos / ingeniería de features

Datos faltantes: Se analizaron los datos faltantes por variable. La mayor parte de datos faltantes se da en las variables '*Sunshine*', '*Evaporation*', '*Cloud3pm*' y '*Cloud9am*'. Mediante agrupamiento de datos por ciudades, se observa que varias presentan datos faltantes en todas sus muestras. Esto podría deberse probablemente a la ausencia del sensor correspondiente en la central de adquisición de datos de esas ciudades.

Imputación de datos faltantes: Se testearon 3 opciones diferentes: supresión de muestras con al menos 1 NaN, imputación por media o moda, supresión de las columnas con mayor cantidad de NaN. Si bien el primer caso descartaría aproximadamente el 40% del dataset, es el que mejor resultado dio. Esto puede deberse a que el número de muestras es muy elevado, las columnas con mayor cantidad de NaN son de vital importancia para el problema y la imputación por media/moda no es adecuada para la cantidad de datos faltantes en este caso.

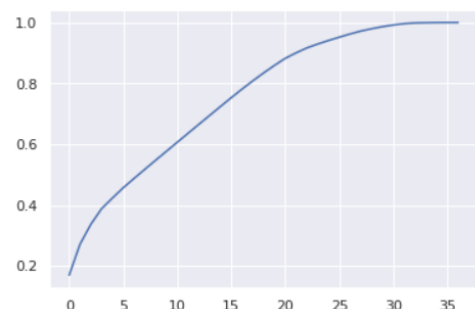
Codificación de la ubicación: Para el caso de la ubicación, se utilizó la API Nominatim de OpenStreetMaps para transformar los nombres de las ciudades en valores de latitud y longitud. Esto no solo reduce drásticamente el número de *features* comparado con la aplicación de OneHotEncoding, sino que también mejora el desempeño ya que la ubicación geográfica permite asociar datos de ciudades cercanas.

Codificación de la dirección del viento: Se transformaron las siglas de dirección en valores de ángulos sexagesimales tomando como referencia 0° el este. A partir de estos valores, se guardan en 2 *features* diferentes los valores de seno y coseno del ángulo. Esto permite al modelo relacionar direcciones de viento similares.

Codificación de la fecha: Se conservó únicamente el mes ya que la estación del año tiene una incidencia fuerte en el clima. Para su codificación se utilizó OneHotEncoding.

Codificación de variables booleanas: Se transformaron en valores numéricos 1/0.

Análisis de features: Se utilizó la matriz de correlación para estudiar la relación entre las variables de entrada y PCA para analizar si existe una jerarquía de features y es posible una reducción de dimensionalidad del dataset. Sin embargo, observando el gráfico acumulado de explicabilidad de la varianza según el número de componentes, vemos que no hay una jerarquía acentuada.



Normalización: Se normalizaron las variables de entrada a media cero y desvío estándar unitario.

4. Entrenamiento de modelos

Se realizaron pruebas con diferentes modelos de clasificación y finalmente se seleccionaron los siguientes 3: Regresión Logística, Random Forest y Gradient Boosting (LGBM), siendo este último el que mejor resultados obtuvo.

Modelo	Accuracy
Regresión Logística	85.43%
Random Forest	86.37%
Gradient Boosting (LGBM)	86.56%

5. Resultados de modelos

Se detallan a continuación los resultados del modelo con mejor desempeño (LGBM Classifier).

	Valor real positivo (Llovió)	Valor real negativo (No llovió)
Predicción positiva (Lloverá)	True Positive 2135	False Positive 692
Predicción negativa (No lloverá)	False Negative 1586	True Negative 12538

Accuracy: 86,56%

Precision: 75,52%

Recall: 57,38%

Utilizando el método `classification_report` de la librería `scikit-learn` se obtiene la siguiente tabla:

```
print(classification_report(y_test, y_pred))
```

	precision	recall	f1-score	support
False	0.89	0.95	0.92	13230
True	0.76	0.57	0.65	3721
accuracy			0.87	16951
macro avg	0.82	0.76	0.78	16951
weighted avg	0.86	0.87	0.86	16951

6. Conclusiones

Se lograron aplicar los conceptos aprendidos durante el cursado de la asignatura a un caso práctico real y se obtuvieron resultados muy buenos.

A medida que se desarrolló el trabajo, se intentaron diferentes métodos de codificación de variables categóricas y de imputación de datos faltantes y se analizó su impacto en el número de features, velocidad de entrenamiento y desempeño de los modelos seleccionados.

Actualmente, el código se encuentra con la configuración de mejor desempeño, pero están comentadas a lo largo del código algunas alternativas que fueron testeadas, pero dieron peores resultados.