



Tecnológico de Monterrey

Análisis de Regresión Lineal Aplicado a un Dataset de Calidad del Aire

Adrián Matute Beltrán A01703889

24 de agosto del 2024

El presente reporte tiene como objetivo analizar las relaciones entre diversas variables medioambientales registradas en un área urbana utilizando un modelo de regresión lineal. El dataset empleado en este análisis, obtenido del UCI Machine Learning Repository, contiene datos recopilados durante un año en una zona urbana con alta contaminación atmosférica. La variable objetivo de este estudio es la concentración de Benceno ($C_6H_6(GT)$), y se busca predecir su comportamiento en función de otras variables atmosféricas y respuestas de sensores químicos.

El propósito de este análisis es determinar la precisión de un modelo de regresión lineal al predecir la concentración de Benceno, identificar las variables más influyentes en su predicción y evaluar la efectividad del modelo.

Descripción del Dataset

El dataset utilizado para el análisis es “Air Quality Data”, contiene 9358 registros utilizables de respuestas promediadas por hora de varios sensores químicos de óxido de metal, junto con datos de un analizador de referencia certificado. Las variables incluidas en el análisis son:

- $CO(GT)$: Concentración horaria promedio de CO en mg/m^3 (medida de referencia).
- $PT08.S1(CO)$: Respuesta del sensor dirigido nominalmente a CO.
- $NMHC(GT)$: Concentración horaria promedio de hidrocarburos no metánicos en $microg/m^3$.
- $C_6H_6(GT)$: Concentración horaria promedio de Benceno en $microg/m^3$.
- $PT08.S2(NMHC)$: Respuesta del sensor dirigido nominalmente a NMHC.
- $NO_x(GT)$: Concentración horaria promedio de NO_x en ppb.
- $PT08.S3(NO_x)$: Respuesta del sensor dirigido nominalmente a NO_x .
- $NO_2(GT)$: Concentración horaria promedio de NO_2 en $microg/m^3$.
- $PT08.S4(NO_2)$: Respuesta del sensor dirigido nominalmente a NO_2 .
- $PT08.S5(O_3)$: Respuesta del sensor dirigido nominalmente a O_3 .
- T: Temperatura en $^{\circ}C$.
- RH: Humedad relativa en %.
- AH: Humedad absoluta.

Se realizó un proceso de preprocesamiento para manejar los valores faltantes y preparar los datos para el análisis.

Preprocesamiento de Datos

El preprocesamiento de los datos fue crucial para asegurar la calidad del análisis. A continuación se describen las etapas de este proceso:

1. Conversión de Tipos de Datos: Algunas columnas contenían datos numéricos representados como cadenas de texto. Se realizó la conversión de éstos valores reemplazando comas por puntos decimales y transformando los datos a tipo 'float64'.

<bound method NDFrame.describe of				Date	Time	CO(GT)	PT08.S1(CO)	NMHC(GT)	C6H6(GT)	PT08.S2(NMHC)	...	PT08.S3(NOx)	NO2(GT)	PT08.S4(NO2)	
PT08.S5(O3)	T	RH	AH												
0	10/03/2004	18.00.00	2.6	1360.0	150.0	11.9	1046.0	...	1056.0	113.0	1692.0	1268.0	13.6	48.9	0.7578
1	10/03/2004	19.00.00	2.0	1292.0	112.0	9.4	955.0	...	1174.0	92.0	1559.0	972.0	13.3	47.7	0.7255
2	10/03/2004	20.00.00	2.2	1402.0	88.0	9.0	939.0	...	1140.0	114.0	1555.0	1074.0	11.9	54.0	0.7502
3	10/03/2004	21.00.00	2.2	1376.0	80.0	9.2	948.0	...	1092.0	122.0	1584.0	1203.0	11.0	60.0	0.7867
4	10/03/2004	22.00.00	1.6	1272.0	51.0	6.5	836.0	...	1205.0	116.0	1490.0	1110.0	11.2	59.6	0.7888
...
1226	30/04/2004	20.00.00	4.4	1449.0	501.0	19.5	1282.0	...	625.0	133.0	2100.0	1569.0	19.1	61.1	1.3345
1227	30/04/2004	21.00.00	3.1	1363.0	234.0	15.1	1152.0	...	684.0	110.0	1951.0	1495.0	18.2	65.4	1.3529
1228	30/04/2004	22.00.00	3.0	1371.0	212.0	14.6	1136.0	...	689.0	102.0	1927.0	1471.0	18.1	66.1	1.3579
1229	30/04/2004	23.00.00	3.1	1406.0	275.0	13.7	1107.0	...	718.0	108.0	1872.0	1384.0	17.7	66.9	1.3422
1230	01/05/2004	00.00.00	3.5	1425.0	275.0	15.2	1155.0	...	709.0	110.0	1936.0	1789.0	17.8	66.8	1.3460

2. Manejo de Valores Faltantes: Los valores faltantes en las columnas de interés, representados por el valor -200, fueron reemplazados por NaN. Posteriormente, se eliminaron todas las filas que contenían estos valores faltantes para asegurar un dataset limpio y completo para el análisis.
3. Selección de Variables: Se seleccionaron las siguientes variables como predictoras: T, RH, AH, PT08.S1(CO), PT08.S2(NMHC), y PT08.S3(NOx). La variable objetivo fue la concentración de Benceno (C6H6(GT)).

Metodología

Para modelar la relación entre las variables predictoras y la concentración de Benceno, se implementó un modelo de regresión lineal desde cero, este enfoque iterativo ajusta los coeficientes del modelo de forma progresiva, minimizando el error cuadrático medio (MSE) en cada iteración.

Implementación del Modelo

1. Normalización de los Datos:
 - a. Las variables predictoras fueron normalizadas para asegurar que todas estén en una escala comparable, lo que ayuda a que el Descenso de Gradiente se haga de manera estable.
2. Inicialización de parámetros
 - a. Los coeficientes del modelo fueron inicializados en cero.
3. Descenso de Gradiente
 - a. Durante 2000 épocas (Epochs), los coeficientes se ajustan utilizando una tasa de aprendizaje de 0.001, para minimizar el MSE.
4. Evaluación del modelo:
 - a. Al finalizar el entrenamiento, se calcula el MSE y el coeficiente de determinación (R^2) para evaluar la precisión del modelo.

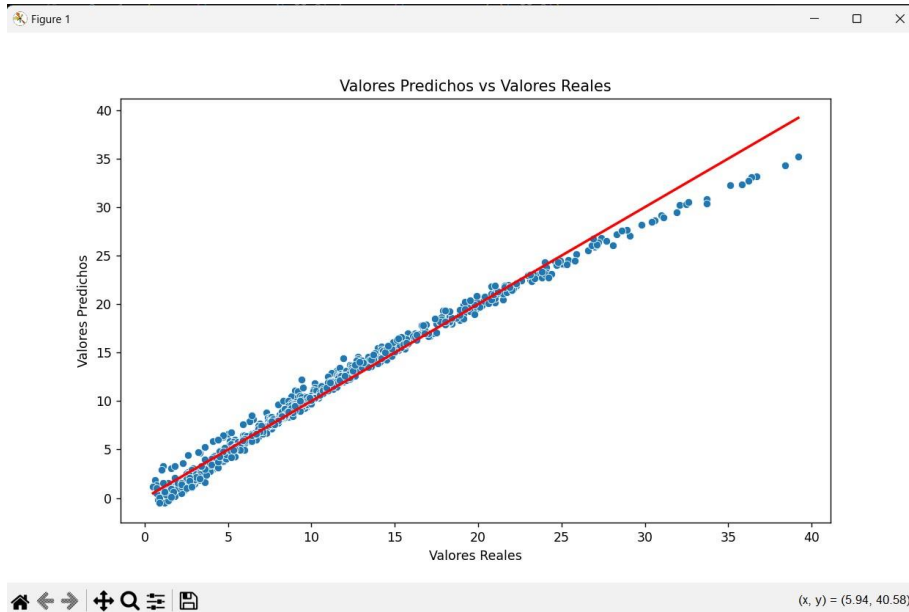
Evaluación del Modelo

El modelo fue evaluado mediante dos métricas principales:

- Error Cuadrático Medio (MSE): Mide el promedio de los errores al cuadrado entre los valores reales y los predichos.
- Coeficiente de Determinación (R^2): Indica la proporción de la varianza en 'Y' explicada por 'X', proporcionando una medida de la calidad del ajuste del modelo

```
Epoch 1100: MSE = 4.791525278329239
Epoch 1200: MSE = 4.23483368610458
Epoch 1300: MSE = 3.844767808992328
Epoch 1400: MSE = 3.5676576747050333
Epoch 1500: MSE = 3.3673109150800897
Epoch 1600: MSE = 3.219296890950137
Epoch 1700: MSE = 3.1071058004638337
Epoch 1800: MSE = 3.019565800172328
Epoch 1900: MSE = 2.9491050716281246
Parámetros finales: [10.57460937 0.07579783 -0.23450221 -0.32265518 2.77571209 3.49916808
-1.10828163]
MSE: 2.8905817628311836
R²: 0.9474077995212299
```

Visualización de los Resultados



Este gráfico muestra la relación entre los valores predichos por el modelo y los valores reales de C₆H₆(GT). Los puntos que se alinean con la línea de identidad (en rojo) indican predicciones precisas.

Interpretación: La alineación cercana de los puntos a la línea de identidad sugiere que el modelo realiza predicciones bastante precisas, aunque pueden observarse algunas desviaciones.

Discusión

El análisis demuestra que la regresión lineal es una herramienta eficaz para predecir la concentración de Benceno en función de otras variables ambientales y de sensores. Sin embargo, algunas desviaciones observadas en el gráfico de predicciones vs valores reales sugieren que podría haber margen para mejorar el modelo.

Referencias

- UCI Machine Learning Repository. Air Quality Data Set. Disponible en:
<https://archive.ics.uci.edu/ml/datasets/air+quality>