

INFORME SOBRE MODELOS DE REGRESIÓN LINEAL

1 DEFINICIÓN MATEMÁTICA

Para un modelo de regresión lineal simple (una sola variable independiente), la relación entre la variable dependiente y y la variable independiente x se describe mediante la siguiente ecuación:

$$y = \beta_0 + \beta_1 x + \epsilon$$

donde:

- y es la variable dependiente.
- x es la variable independiente.
- β_0 es la intersección o término constante (intercepto), que representa el valor esperado de y cuando $x = 0$.
- β_1 es la pendiente de la línea de regresión, que representa el cambio esperado en y por unidad de cambio en x .
- ϵ es el término de error aleatorio, que captura la desviación de los puntos de datos reales de la línea de regresión debido a factores no observados.

1.1 Suposiciones del Modelo

1. **Linealidad:** La relación entre la variable dependiente y y la variable independiente x es lineal.
2. **Independencia:** Las observaciones son independientes entre sí.
3. **Homoscedasticidad:** La varianza de los errores es constante para todos los valores de x .
4. **Normalidad:** Los errores (ϵ) están normalmente distribuidos con media cero.

1.2 Estimación de los Parámetros

Para estimar los parámetros β_0 y β_1 , se utiliza el método de los mínimos cuadrados ordinarios (OLS), que minimiza la suma de los cuadrados de los errores (residuos). La fórmula para los estimadores de mínimos cuadrados de β_0 y β_1 son:

$$\hat{\beta}_1 = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2}$$

$$\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x}$$

donde:

- $\hat{\beta}_0$ es el estimador de β_0 .
- $\hat{\beta}_1$ es el estimador de β_1 .
- \bar{x} es la media de los valores de x .
- \bar{y} es la media de los valores de y .
- n es el número de observaciones.

1.3 Modelo de Regresión Lineal Múltiple

Para un modelo de regresión lineal múltiple (varias variables independientes), la relación entre la variable dependiente y y las variables independientes x_1, x_2, \dots, x_n se describe mediante la siguiente ecuación:

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_n x_n + \epsilon$$

donde:

- y es la variable dependiente.
- x_1, x_2, \dots, x_n son las variables independientes.
- β_0 es la intersección o término constante (intercepto).
- $\beta_1, \beta_2, \dots, \beta_n$ son los coeficientes de regresión para las variables independientes.
- ϵ es el término de error aleatorio.

1.4 Estimación de los Parámetros en el Modelo Múltiple

En el modelo de regresión lineal múltiple, los parámetros $\beta_0, \beta_1, \dots, \beta_n$ se estiman utilizando también el método de los mínimos cuadrados ordinarios (OLS), que minimiza la suma de los cuadrados de los errores. Las fórmulas para los estimadores en el modelo múltiple son más complejas y requieren el uso de álgebra matricial.

1.5 Evaluación del Modelo

Una vez estimados los parámetros, es importante evaluar el modelo utilizando métricas como:

- **Coefficiente de determinación (R^2):** Indica la proporción de la variabilidad en la variable dependiente explicada por las variables independientes.
- **Análisis de residuos:** Para verificar las suposiciones de normalidad, homoscedasticidad e independencia.
- **Pruebas de significancia estadística:** Para determinar si los coeficientes de regresión son significativamente diferentes de cero.

2 DESCRIPCIÓN DEL DATASET

2.1 Fallecidos por COVID-19 - [Ministerio de Salud - MINSA]

Es el registro diario de muertes por Covid-19. Cada registro es igual a una persona, la cual puede caracterizarse por sexo, edad y ubicación geográfica hasta nivel de distrito; además, el 06.mayo.2021 se agregó el código UBIGEO.

Desde que se publicó este dataset, cada registro representaba un fallecido confirmado por covid-19, quienes cumplen con criterios clínicos y de laboratorio (prueba molecular, antigénica o pruebas serológicas). A partir del 31.mayo.2021 se cambió el criterio de “Fallecidos por Covid-19” por “Muertes por Covid-19” y como resultado el dataset creció casi al triple en el número de registros. Esta nueva clasificación está definida por el cumplimiento de al menos uno de los siguientes siete criterios técnicos:

- **Criterio virológico:** Muerte en un caso confirmado de COVID-19 que fallece en los 60 días posteriores a una prueba molecular (PCR) o antigénica reactiva para SARS-CoV-2.
- **Criterio serológico:** Muerte en un caso confirmado de COVID-19 que fallece en los 60 días posteriores a una prueba serológica positiva IgM o IgM/IgG para SARS-CoV-2.

- **Criterio radiológico:** Muerte en un caso probable de COVID-19 que presenta una imagen radiológica, tomográfica o de resonancia magnética nuclear compatible con neumonía COVID-19.
- **Criterio nexa epidemiológico:** Muerte en un caso probable de COVID-19 que presenta nexa epidemiológico con un caso confirmado de COVID-19.
- **Criterio investigación epidemiológica:** Muerte en un caso sospechoso de COVID-19 que es verificado por investigación epidemiológica de la Red Nacional de Epidemiología (RENACE).
- **Criterio clínico:** Muerte en un caso sospechoso de COVID-19 que presenta cuadro clínico compatible con la enfermedad.
- **Criterio SINADEF:** Muerte con certificado de defunción en el que se presenta el diagnóstico de COVID-19 como causa de la muerte. El fallecimiento por COVID-19 en el certificado de defunción está definido por la presencia en los campos A, B, C o D de los códigos CIE-10: U071, U072, B342, B972, o la mención de los términos “coronavirus”, “cov-2”, “cov2”, “covid” y “sars”.

3 CÓDIGO

```

1 import streamlit as st
2 import pandas as pd
3 import matplotlib.pyplot as plt
4 from sklearn.linear_model import LinearRegression
5 from sklearn.metrics import mean_squared_error
6 import pdfkit
7
8 def main():
9     st.title("Regresi n Lineal con Streamlit")
10
11     # Cargar el dataset
12     uploaded_file = st.file_uploader("Carga tu dataset", type=['csv'])
13     if uploaded_file is not None:
14         df = pd.read_csv(uploaded_file)
15         st.write("Dataset cargado:")
16         st.write(df.head())
17
18         # Seleccionar las columnas para la regresión
19         numeric_columns = df.select_dtypes(include=['int64', 'float64']).columns
20         x_col = st.selectbox("Selecciona la variable independiente (X)", numeric_columns)
21         y_col = st.selectbox("Selecciona la variable dependiente (Y)", numeric_columns)
22
23         if st.button("Realizar regresión lineal"):
24             # Entrenar el modelo de regresión lineal
25             X = df[x_col].values.reshape(-1, 1)
26             y = df[y_col].values
27             model = LinearRegression()
28             model.fit(X, y)
29
30             # Mostrar los resultados
31             st.write(f"Coeficiente de regresión (m): {model.coef_[0]:.2f}")
32             st.write(f"Intercepto independiente (b): {model.intercept_:.2f}")
33             st.write(f"R-squared: {model.score(X, y):.2f}")
34
35             # Interpretación estadística
36             st.write("Interpretación estadística:")
37             st.write(f"El coeficiente de regresión (m) indica que por cada unidad de aumento en {x_col}, {y_col} aumenta en {model.coef_[0]:.2f} unidades.")
38             st.write(f"El intercepto independiente (b) indica que cuando {x_col} es 0, {y_col} es igual a {model.intercept_:.2f}.")
39             st.write(f"El R-squared indica que el {model.score(X, y)*100:.2f}% de la variabilidad en {y_col} se explica por la variabilidad en {x_col}.")
40
41             # Graficar los resultados
42             fig, ax = plt.subplots(figsize=(8, 6))

```

```

43     ax.scatter(X, y, label='Datos')
44     ax.plot(X, model.predict(X), color='red', label='Regresi n Lineal')
45     ax.set_xlabel(x_col)
46     ax.set_ylabel(y_col)
47     ax.set_title("Regresi n Lineal")
48     ax.legend()
49     st.pyplot(fig)
50
51     # Proyecciones
52     y_pred = model.predict(X)
53     fig, ax = plt.subplots(figsize=(8, 6))
54     ax.plot(y, label='Valor real')
55     ax.plot(y_pred, label='Proyecci n')
56     ax.set_xlabel(' ndice ')
57     ax.set_ylabel(y_col)
58     ax.set_title("Proyecciones")
59     ax.legend()
60     st.pyplot(fig)
61
62     # M tricas de evaluaci n
63     mse = mean_squared_error(y, y_pred)
64     st.write(f"Mean Squared Error: {mse:.2f}")
65     st.write(f"La ra z cuadrada del error cuadr tico medio es {mse**0.5:.2f}, lo
que indica que el modelo tiene un error de aproximadamente {mse**0.5:.2f} unidades en la
predicci n de {y_col}.")
66
67
68 if __name__ == "__main__":
69     main()

```

Listing 1: Código Python

4 RESULTADOS

Dataset cargado:

	FECHA_CORTE	FECHA_FALLECIMIENTO	EDAD_DECLARADA	SEXO	CLASIFICACION_DEF
0	20,240,317	20,220,219	63	MASCULINO	Criterio virológico
1	20,240,317	20,210,529	74	MASCULINO	Criterio virológico
2	20,240,317	20,210,623	72	FEMENINO	Criterio SINADef
3	20,240,317	20,210,824	85	MASCULINO	Criterio investigación Epidem
4	20,240,317	20,210,627	46	MASCULINO	Criterio virológico

4.1 Interpretacion

Selecciona la variable independiente (X)

FECHA_FALLECIMIENTO

Selecciona la variable dependiente (Y)

EDAD_DECLARADA

Realizar regresión lineal

Coefficiente de regresión (m): 0.00

Término independiente (b): -1564.98

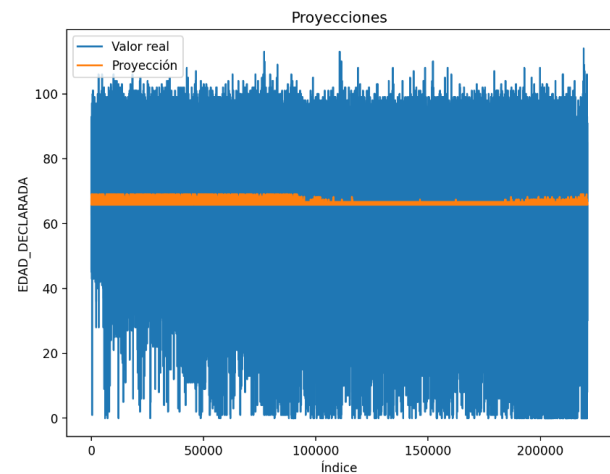
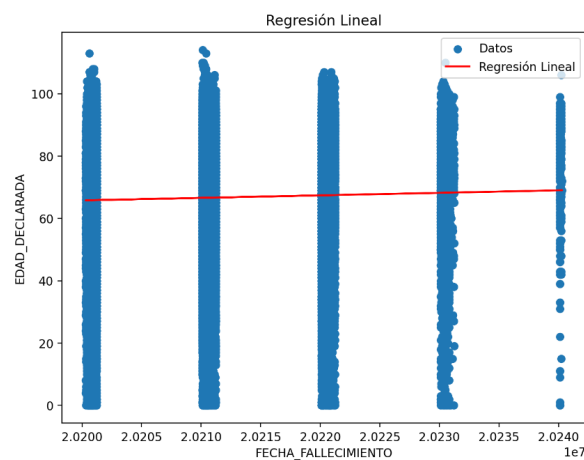
R-squared: 0.00

Interpretación estadística:

El coeficiente de regresión (m) indica que por cada unidad de aumento en FECHA_FALLECIMIENTO, EDAD_DECLARADA aumenta en 0.00 unidades.

El término independiente (b) indica que cuando FECHA_FALLECIMIENTO es 0, EDAD_DECLARADA es igual a -1564.98.

El R-squared indica que el 0.11% de la variabilidad en EDAD_DECLARADA se explica por la variabilidad en FECHA_FALLECIMIENTO.



5 MÉTRICA DE EVALUACIÓN

En la regresión lineal para evaluarlo se utilizan métricas como el Error Cuadrático Medio (MSE), el Coeficiente de Determinación (R^2), entre otros, que miden qué tan bien el modelo ajusta los datos reales con respecto a las predicciones numéricas.

Mean Squared Error: 251.87

La raíz cuadrada del error cuadrático medio es 15.87, lo que indica que el modelo tiene un error de aproximadamente 15.87 unidades en la predicción de EDAD_DECLARADA.

Este trabajo se desarrollo con la aplicación de la Streamlit.