

Modelos y técnicas de Minería de Datos para Regresión

Profesor: Juan Gamarra Moreno

Definición de Regresión

La regresión en el contexto de la minería de datos es una técnica estadística y de aprendizaje automático utilizada para modelar y analizar la relación entre una variable dependiente continua (también llamada variable de respuesta) y una o más variables independientes (también conocidas como predictores o características). El objetivo principal de la regresión es predecir el valor de la variable dependiente basándose en los valores de las variables independientes.

Importancia de la Regresión (1/2)

- 1. Predicción y Pronóstico:** La regresión es fundamental para predecir valores futuros de una variable dependiente basándose en valores presentes y pasados de las variables independientes. Esto es crucial en campos como la economía, las finanzas y la meteorología, donde las predicciones precisas son esenciales para la toma de decisiones.
- 2. Análisis de Relaciones:** Permite analizar y cuantificar la relación entre variables. Esto ayuda a entender cómo una variable independiente puede afectar a una variable dependiente, lo cual es útil en investigaciones científicas y estudios de mercado.
- 3. Modelado de Datos Complejos:** La regresión, especialmente en sus formas avanzadas como la regresión polinómica y las técnicas de regularización, permite capturar relaciones complejas y no lineales entre las variables, ofreciendo modelos más precisos y robustos.

Importancia de la Regresión (2/2)

- 4. Identificación de Factores Relevantes:** Ayuda a identificar y seleccionar las variables más importantes que influyen en la variable dependiente, facilitando la simplificación de modelos y la eliminación de variables irrelevantes.
- 5. Evaluación y Mejora de Procesos:** En la industria y la manufactura, la regresión se utiliza para evaluar y mejorar procesos, optimizando la eficiencia y reduciendo costos mediante el análisis de factores que afectan la producción y la calidad.

Aplicaciones de la Regresión

1. Economía y Finanzas:

- **Predicción de Precios:** Utilizada para predecir precios de acciones, bienes raíces y otros activos financieros.
- **Análisis de Riesgos:** Modela el riesgo financiero y evalúa la relación entre diferentes factores económicos.

2. Marketing y Ventas:

- **Pronóstico de Ventas:** Estima futuras ventas basándose en datos históricos y tendencias actuales.
- **Análisis de Clientes:** Identifica los factores que influyen en el comportamiento de compra de los clientes, ayudando en la segmentación de mercados y en la personalización de estrategias de marketing.

Aplicaciones de la Regresión

3. Salud y Medicina:

- **Modelado de Resultados Clínicos:** Predice resultados de tratamientos médicos basándose en variables como edad, historial médico y tratamientos anteriores.
- **Epidemiología:** Analiza la relación entre factores de riesgo y la incidencia de enfermedades, ayudando a formular políticas de salud pública.

4. Ingeniería y Manufactura:

- **Control de Calidad:** Modela la relación entre variables de proceso y la calidad del producto final, ayudando a identificar y corregir problemas en la producción.
- **Optimización de Procesos:** Utiliza la regresión para mejorar la eficiencia de procesos industriales y de manufactura.

Aplicaciones de la Regresión

5. Ciencias Sociales:

- **Investigación Educativa:** Analiza el impacto de variables como el nivel socioeconómico, el tipo de escuela y el entorno familiar en el rendimiento académico de los estudiantes.
- **Sociología:** Estudia las relaciones entre variables sociales como el ingreso, la educación y la salud.

6. Ciencias Ambientales:

- **Predicción del Clima:** Modela y predice variables climáticas basándose en datos históricos y patrones actuales.
- **Análisis de Impacto Ambiental:** Evalúa cómo las actividades humanas influyen en variables ambientales como la calidad del aire y el agua.

Aplicaciones de la Regresión

5. Ciencias Sociales:

- **Investigación Educativa:** Analiza el impacto de variables como el nivel socioeconómico, el tipo de escuela y el entorno familiar en el rendimiento académico de los estudiantes.
- **Sociología:** Estudia las relaciones entre variables sociales como el ingreso, la educación y la salud.

6. Ciencias Ambientales:

- **Predicción del Clima:** Modela y predice variables climáticas basándose en datos históricos y patrones actuales.
- **Análisis de Impacto Ambiental:** Evalúa cómo las actividades humanas influyen en variables ambientales como la calidad del aire y el agua.

Aplicaciones de la Regresión

8. Deportes:

- **Análisis de Rendimiento:** Predice el rendimiento de atletas basándose en variables como el entrenamiento, la nutrición y el historial de lesiones.
- **Estrategias de Juego:** Modela y analiza estrategias de juego, optimizando decisiones tácticas y estratégicas.

La regresión es una herramienta versátil y poderosa en la minería de datos que se aplica en una amplia gama de disciplinas para la toma de decisiones informadas, la optimización de procesos y la mejora de resultados

Fundamentos Teóricos de la Regresión

Conceptos básicos (1/4)

Variables Dependientes

La variable dependiente, también conocida como la variable de respuesta, es la variable que se intenta predecir o explicar. Su valor depende de las variables independientes. En el contexto de la regresión, es la variable que se coloca en el eje *y* de un gráfico.

Ejemplos:

- En un estudio sobre el impacto de la educación en el ingreso, el ingreso sería la variable dependiente.
- En la predicción del precio de una casa, el precio sería la variable dependiente.

Conceptos básicos (2/4)

Variables Independientes

Las variables independientes, también conocidas como predictores o explicativas, son las variables que se utilizan para predecir o explicar la variable dependiente. Estas variables se colocan en el eje x de un gráfico.

Ejemplos:

- En el mismo estudio sobre educación e ingreso, los años de educación, la experiencia laboral y la ocupación serían variables independientes.
- En la predicción del precio de una casa, características como el tamaño de la casa, el número de habitaciones y la ubicación serían variables independientes.

Conceptos básicos (3/4)

Relación Lineal

Una relación lineal implica que hay una relación directa y proporcional entre las variables independientes y la variable dependiente. Esto significa que un cambio en la variable independiente resulta en un cambio constante en la variable dependiente. La relación puede representarse con una línea recta en un gráfico.

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

Donde:

- y es la variable dependiente.
- β_0 es el intercepto.
- β_1 es la pendiente de la línea.
- x es la variable independiente.
- ϵ es el término de error.

En este modelo, si x aumenta en una unidad, y aumenta en β_1 unidades.

Conceptos básicos (4/4)

Relación No Lineal

Una relación no lineal implica que el cambio en la variable dependiente no es proporcional al cambio en las variables independientes. Estas relaciones pueden adoptar diversas formas, como curvas, parábolas, exponenciales, etc. Las relaciones no lineales no se pueden representar adecuadamente con una línea recta.

Ejemplo: $y = \beta_0 + \beta_1 x + \beta_2 x^2 + \epsilon$

En este modelo cuadrático:

- y es la variable dependiente.
- β_0 es el intercepto.
- β_1 es el coeficiente del término lineal.
- β_2 es el coeficiente del término cuadrático.
- x es la variable independiente.
- ϵ es el término de error.

Aquí, la relación entre x e y es curva. A medida que x cambia, el cambio en y depende tanto del valor de x como de x^2 .

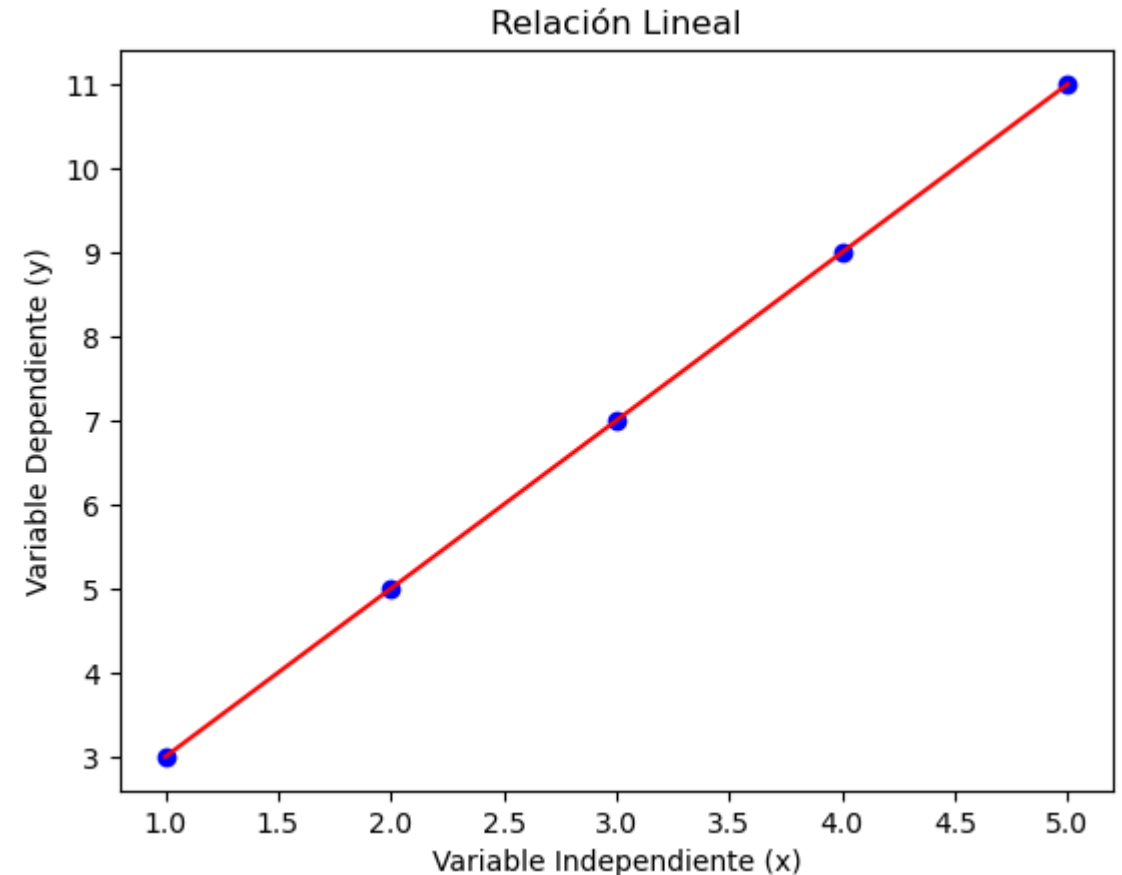
Visualización de Relaciones

- **Relación Lineal:** En un gráfico de dispersión, una relación lineal entre x e y se verá como una línea recta.
- **Relación No Lineal:** En un gráfico de dispersión, una relación no lineal entre x e y se verá como una curva.

```
import numpy as np
import matplotlib.pyplot as plt

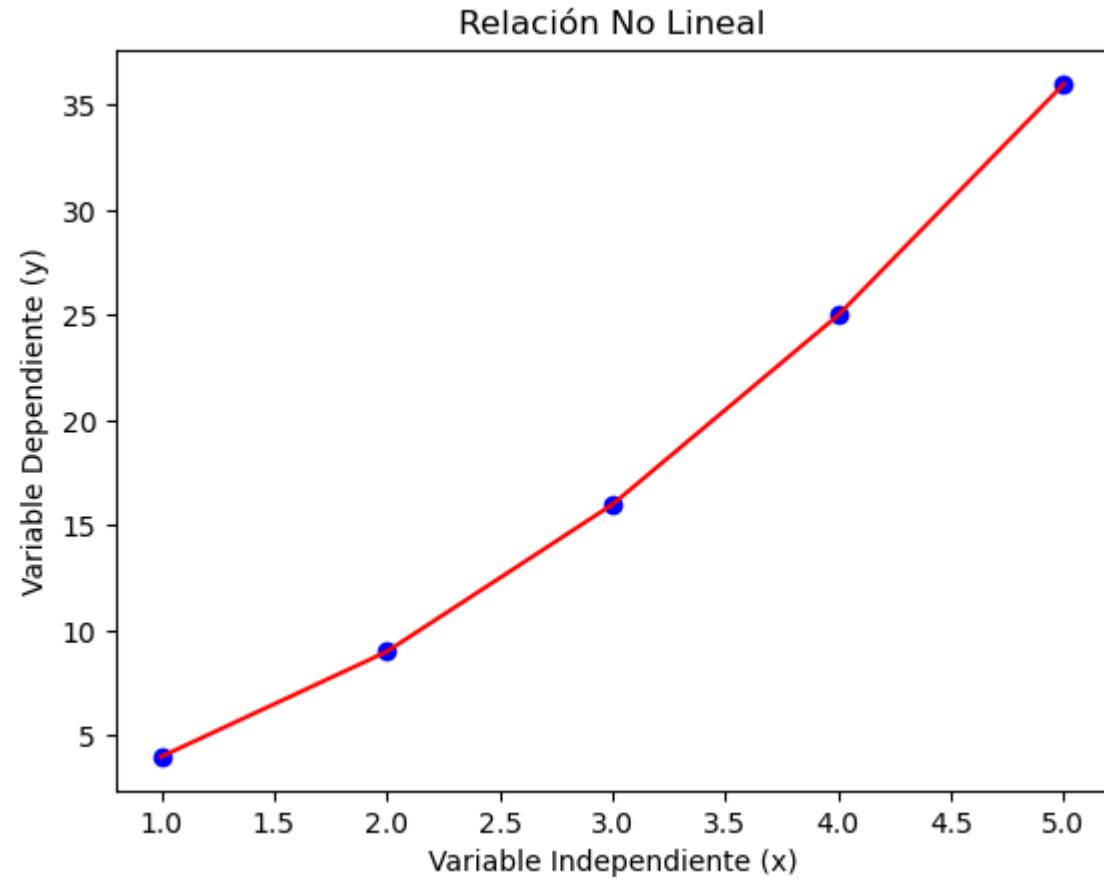
# Datos de ejemplo
x = np.array([1, 2, 3, 4, 5])
y = 2 * x + 1

plt.scatter(x, y, color='blue')
plt.plot(x, y, color='red')
plt.xlabel('Variable Independiente (x)')
plt.ylabel('Variable Dependiente (y)')
plt.title('Relación Lineal')
plt.show()
```




```
# Datos de ejemplo
x = np.array([1, 2, 3, 4, 5])
y = x**2 + 2 * x + 1

plt.scatter(x, y, color='blue')
plt.plot(x, y, color='red')
plt.xlabel('Variable Independiente (x)')
plt.ylabel('Variable Dependiente (y)')
plt.title('Relación No Lineal')
plt.show()
```



Importante

En resumen, entender la distinción entre variables dependientes e independientes y las diferencias entre relaciones lineales y no lineales es crucial para el análisis de datos y la construcción de modelos predictivos precisos en la minería de datos.

Modelos de Regresión Lineal

Regresión Lineal Simple

La regresión lineal simple es un método estadístico que explora la relación entre dos variables: una variable independiente x y una variable dependiente y . El objetivo es ajustar una línea recta que minimice la suma de los cuadrados de las diferencias entre los valores observados y los valores predichos por el modelo.

Ecuación del Modelo: $y = \beta_0 + \beta_1 x + \epsilon$

Donde:

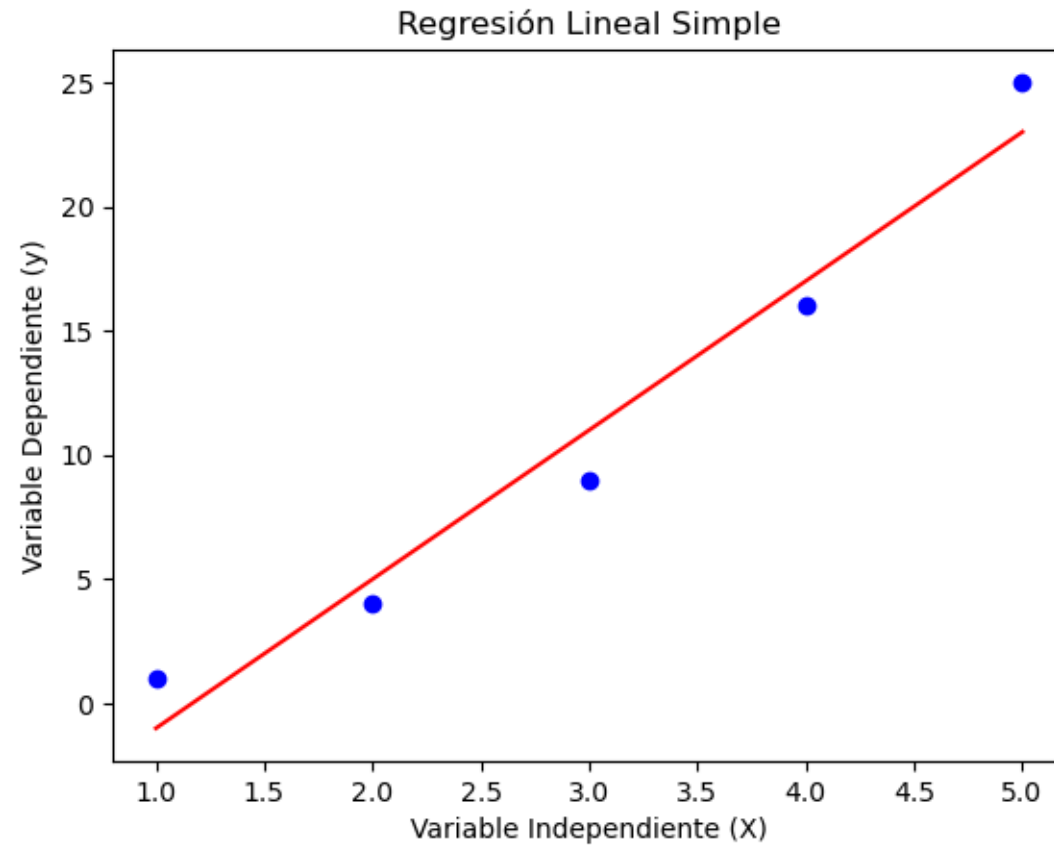
- y es la variable dependiente.
- β_0 es el intercepto.
- β_1 es la pendiente de la línea.
- x es la variable independiente.
- ϵ es el término de error.

Regresión Lineal Simple

Supuestos:

1. **Linealidad:** La relación entre x e y es lineal.
2. **Independencia:** Las observaciones son independientes entre sí.
3. **Homocedasticidad:** La varianza de los errores es constante a lo largo de las observaciones.
4. **Normalidad de los Errores:** Los errores siguen una distribución normal.

Regresión Lineal Simple



Regresión Lineal Múltiple

La regresión lineal múltiple extiende la regresión lineal simple al usar múltiples variables independientes para predecir una variable dependiente. Este método es útil cuando se quiere entender cómo varias variables explicativas influyen en una variable de respuesta.

Ecuación del Modelo: $y = \beta_0 + \beta_1x_1 + \beta_2x_2 + \dots + \beta_nx_n + \epsilon$

Donde:

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

Regresión Lineal Múltiple

Supuestos:

1. **Linealidad:** La relación entre las variables independientes y la variable dependiente es lineal.
2. **Independencia:** Las observaciones son independientes entre sí.
3. **Homocedasticidad:** La varianza de los errores es constante a lo largo de las observaciones.
4. **Normalidad de los Errores:** Los errores siguen una distribución normal.
5. **No Multicolinealidad:** Las variables independientes no están altamente correlacionadas entre sí.

Regresión Lineal Múltiple

Comparación entre Regresión Lineal Simple y Múltiple

Característica	Regresión Lineal Simple	Regresión Lineal Múltiple
Número de variables independientes	Una sola variable independiente	Dos o más variables independientes
Complejidad	Menos compleja, más fácil de interpretar	Más compleja, pero puede capturar más factores explicativos
Aplicaciones	Cuando se conoce un solo factor principal que influye en y	Cuando hay múltiples factores que influyen en y
Supuestos adicionales	Ninguno específico	No multicolinealidad entre las variables independientes

En resumen, la elección entre la regresión lineal simple y múltiple depende del contexto y de la cantidad de variables independientes que se consideren relevantes para explicar la variable dependiente. Ambos modelos son fundamentales en la minería de datos y proporcionan una base sólida para el análisis predictivo y el entendimiento de las relaciones entre variables.

Algoritmos de Regresión Lineal

Para los modelos de regresión lineal, se utilizan varios algoritmos y técnicas que facilitan la estimación de los coeficientes y aseguran que el modelo se ajuste adecuadamente a los datos.

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

Donde:

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

Algoritmos de Regresión Lineal

Método de Mínimos Cuadrados Ordinarios (OLS)

Descripción: El método de mínimos cuadrados ordinarios es la técnica más común para estimar los coeficientes de un modelo de regresión lineal. OLS minimiza la suma de los cuadrados de las diferencias entre los valores observados y los valores predichos por el modelo.

Ecuación: $\hat{\beta} = (X^T X)^{-1} X^T y$

Donde:

- $\hat{\beta}$ es el vector de coeficientes estimados.
- X es la matriz de variables independientes.
- y es el vector de la variable dependiente.

Algoritmos de Regresión Lineal

Regresión Ridge

Descripción: La regresión Ridge es una variante de la regresión lineal que incluye una penalización L2 sobre los coeficientes. Este método es útil para manejar la multicolinealidad y prevenir el sobreajuste.

Ecuación: $\hat{\beta} = (X^T X + \lambda I)^{-1} X^T y$

Donde:

- λ es el parámetro de regularización que controla la penalización.
- I es la matriz identidad.

Algoritmos de Regresión Lineal

Regresión Lasso

Descripción: La regresión Lasso (Least Absolute Shrinkage and Selection Operator) es otra variante de la regresión lineal que aplica una penalización L1. Este método no solo previene el sobreajuste, sino que también puede llevar a la eliminación de algunos coeficientes, resultando en un modelo más interpretativo.

Ecuación: Minimizar $\sum_{i=1}^n (y_i - \hat{y}_i)^2 + \lambda \sum_{j=1}^p |\beta_j|$

Donde:

- λ es el parámetro de regularización que controla la penalización.

Algoritmos de Regresión Lineal

Regresión Elastic Net

Descripción: Elastic Net combina las penalizaciones L1 y L2 de Lasso y Ridge. Es útil cuando hay múltiples características correlacionadas.

Ecuación: Minimizar $\sum_{i=1}^n (y_i - \hat{y}_i)^2 + \lambda_1 \sum_{j=1}^p |\beta_j| + \lambda_2 \sum_{j=1}^p \beta_j^2$

Donde:

- λ_1 y λ_2 son los parámetros de regularización.

Algoritmos de Regresión Lineal

Descenso del Gradiente

Descripción: El descenso del gradiente es un algoritmo iterativo que optimiza los coeficientes de un modelo de regresión lineal minimizando una función de pérdida (generalmente el error cuadrático medio). Es especialmente útil para grandes conjuntos de datos.

Ecuación:
$$\beta_j \leftarrow \beta_j - \alpha \frac{\partial J}{\partial \beta_j}$$

Donde:

- α es la tasa de aprendizaje.
- J es la función de costo (generalmente el MSE).

Otros modelos de Regresión

Regresión Polinómica y Características Polinomiales

Regresión Polinómica

- **Definición:** La regresión polinómica es una forma de regresión lineal en la que la relación entre la variable dependiente y y la variable independiente x se modela como un polinomio de grado n . Esto permite capturar relaciones no lineales entre las variables.
- **Ecuación del Modelo:**

Para una regresión polinómica de grado 2, la ecuación es:

$$y = \beta_0 + \beta_1 x + \beta_2 x^2 + \epsilon$$

Para un polinomio de grado n , la ecuación general es:

$$y = \beta_0 + \beta_1 x + \beta_2 x^2 + \dots + \beta_n x^n + \epsilon$$

Regresión Polinómica y Características Polinomiales

Transformación de Variables

- **Concepto:** En la regresión polinómica, las variables independientes originales se transforman en características polinómicas. Esto se hace elevando las variables independientes a diversas potencias y, en algunos casos, creando productos cruzados de variables.
- **Ejemplo de Transformación:** Para un modelo de grado 2 con una variable independiente x , las nuevas características serían x y x^2 .

Regresión Polinómica y Características Polinomiales

Comparación con la Regresión Lineal Simple

- **Regresión Lineal Simple:**

- Asume una relación lineal directa entre la variable independiente x y la variable dependiente y .
- Ecuación del modelo:

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

- **Regresión Polinómica:**

- Permite modelar relaciones no lineales al incluir términos polinómicos de las variables independientes.
- Ecuación del modelo para un polinomio de grado n :

$$y = \beta_0 + \beta_1 x + \beta_2 x^2 + \dots + \beta_n x^n + \epsilon$$

Regresión Polinómica y Características Polinomiales

Ventajas y Desventajas:

Característica	Regresión Lineal Simple	Regresión Polinómica
Relación entre variables	Lineal	No lineal
Complejidad del modelo	Baja	Moderada a alta (dependiendo del grado)
Facilidad de interpretación	Alta	Moderada (mayor complejidad interpretativa)
Flexibilidad	Baja	Alta (mejor ajuste para relaciones complejas)
Riesgo de sobreajuste	Bajo	Alto (especialmente con grados altos)

Modelos de Regresión No Lineal

La regresión no lineal es una técnica de modelado en la que la relación entre la variable dependiente y las variables independientes se modela mediante una función no lineal. A diferencia de la regresión lineal, donde se asume que las variables tienen una relación lineal, en la regresión no lineal, la relación puede ser de cualquier forma que no sea una línea recta.

Importancia:

- **Captura de Relaciones Complejas:** Permite capturar relaciones complejas entre variables que no pueden ser representadas adecuadamente por modelos lineales.
- **Flexibilidad:** Ofrece una mayor flexibilidad en el modelado de datos, especialmente cuando se espera que los datos sigan un patrón curvo o no lineal.

Modelos de Regresión No Lineal

La regresión no lineal es una técnica de modelado en la que la relación entre la variable dependiente y las variables independientes se modela mediante una función no lineal. A diferencia de la regresión lineal, donde se asume que las variables tienen una relación lineal, en la regresión no lineal, la relación puede ser de cualquier forma que no sea una línea recta.

Importancia:

- **Captura de Relaciones Complejas:** Permite capturar relaciones complejas entre variables que no pueden ser representadas adecuadamente por modelos lineales.
- **Flexibilidad:** Ofrece una mayor flexibilidad en el modelado de datos, especialmente cuando se espera que los datos sigan un patrón curvo o no lineal.

Técnicas Avanzadas de Regresión

1. Regresión Ridge y Lasso

- Regularización de modelos

- Comparación de Ridge y Lasso

2. Regresión Elastic Net

- Combinación de Ridge y Lasso

- Ventajas y aplicaciones

Técnicas de Regresión Basadas en Árboles

1. Árboles de decisión para regresión

- Construcción de árboles
- Poda de árboles

2. Random Forest para regresión

- Ensamblaje de árboles
- Importancia de variables

3. Boosting para regresión

- AdaBoost
- Gradient Boosting Machines (GBM)
- XGBoost

Modelos de Regresión Basados en Máquinas de Soporte Vectorial (SVM)

1. Introducción a SVM para regresión (SVR)
2. Función de pérdida Hinge
3. Kernel trick para modelos no lineales

Redes Neuronales para Regresión

1. Introducción a redes neuronales
2. Arquitectura de una red neuronal para regresión
 - Capas y neuronas
 - Función de activación
3. Entrenamiento de redes neuronales
 - Algoritmo de backpropagation
 - Técnicas de optimización

Evaluación y Validación de Modelos de Regresión

1. Validación cruzada
2. Curva de aprendizaje
3. Técnicas de resampling

Tendencias y Avances en Técnicas de Regresión

Regresión cuántica

Regresión con datos de alta dimensionalidad

Aplicaciones en Big Data y aprendizaje automático