

Regresión lineal: el enfoque estadístico con máxima verosimilitud

Matías Carrasco

4 de septiembre de 2023

Índice

1. Modelo lineal gaussiano	1
2. Verosimilitud	2
3. Verosimilitud del conjunto de datos	3
4. Máxima verosimilitud	4
5. Derivada respecto a θ	4
6. Derivada respecto a σ	5
7. Modelos más complejos y sobreajuste	6
8. Incertidumbre en las estimaciones de los parámetros	6
9. Variabilidad en las predicciones	8
10. Varianza de y_{nuevo} en el caso univariado	9

1. Modelo lineal gaussiano

El i -ésimo target y_i se genera de la forma $\theta^\top \mathbf{x}_i$ más una cantidad aleatoria ϵ_i :

$$y_i = \theta^\top \mathbf{x}_i + \epsilon_i$$

Recordar que la ecuación $y_i = \boldsymbol{\theta}^\top \mathbf{x}_i + \epsilon_i$ en forma matricial con los vectores $\boldsymbol{\theta}$ y \mathbf{x}_i dados por

$$\boldsymbol{\theta} = \begin{bmatrix} b \\ w_1 \\ \vdots \\ w_D \end{bmatrix} \quad \mathbf{x}_i = \begin{bmatrix} 1 \\ x_i^{(1)} \\ \vdots \\ x_i^{(D)} \end{bmatrix}$$

representa

$$y_i = [b \quad w_1 \quad \dots \quad w_D] \begin{bmatrix} 1 \\ x_i^{(1)} \\ \vdots \\ x_i^{(D)} \end{bmatrix} + \epsilon_i = b + \sum_{j=1}^D w_j x_i^{(j)} + \epsilon_i$$

La diferencia entre el modelo y el valor real del target, la ϵ_i , es una variable aleatoria continua. Las asumimos independientes entre sí:

$$p(\epsilon_1, \dots, \epsilon_i) = \prod_{i=1}^N p(\epsilon_i)$$

Para la forma de $p(\epsilon_i)$ asumiremos que se trata de una distribución gaussiana (o normal) con media cero y varianza σ^2 . Esto permite que ϵ_i sea tanto positivo como negativo (permite que los datos se encuentren tanto por encima como por debajo de la línea $\boldsymbol{\theta}^\top \mathbf{x}$) y tiene interesantes propiedades de modelado que lo vinculan a la función de pérdida cuadrática.

Nuestro modelo consta entonces de dos componentes:

1. Un componente determinístico $\boldsymbol{\theta}^\top \mathbf{x}_i$.
2. Un componente aleatorio ϵ_i , a veces referido como ruido.

Existen alternativas al ruido gaussiano y también al ruido aditivo. Por ejemplo, para algunas aplicaciones, un término multiplicativo podría ser más apropiado (en cuyo caso, $y = f(\mathbf{x}; \boldsymbol{\theta})\epsilon$). Sin embargo, el ruido aditivo gaussiano nos permite obtener expresiones exactas para el valor óptimo del parámetro $\hat{\boldsymbol{\theta}}$.

2. Verosimilitud

Nuestro modelo tiene la siguiente forma:

$$y_i = f(\mathbf{x}_i; \boldsymbol{\theta}) + \epsilon_i, \quad \epsilon_i \sim \mathcal{N}(0, \sigma^2)$$

Necesitamos encontrar el valor óptimo $\hat{\boldsymbol{\theta}}$ de $\boldsymbol{\theta}$. También tenemos un parámetro adicional σ^2 .

La salida del modelo, y , es ahora en sí misma una variable aleatoria. En particular, no hay un único valor de y_i para un valor particular de \mathbf{x}_i . La variable aleatoria y_i tiene la siguiente función de densidad:

$$p(y_i | \mathbf{x}_i, \boldsymbol{\theta}, \sigma^2) = \mathcal{N}(\boldsymbol{\theta}^\top \mathbf{x}_i, \sigma^2)$$

Notar la condicional en el lado izquierdo - la densidad de y_i depende de valores particulares de \mathbf{x}_i y $\boldsymbol{\theta}$ (determinan la media) y σ^2 (la varianza).

3. Verosimilitud del conjunto de datos

Si tenemos N datos, estamos interesados en la densidad conjunta condicional:

$$p(y_1, \dots, y_N | \mathbf{x}_1, \dots, \mathbf{x}_N, \boldsymbol{\theta}, \sigma^2)$$

Esta es una densidad conjunta sobre todos los targets en nuestro conjunto de datos. Lo escribiremos de manera compacta (usando notación vectorial y la matriz de diseño \mathbf{x}) como

$$p(\mathbf{y} | \mathbf{x}, \boldsymbol{\theta}, \sigma^2).$$

Evaluar esta densidad en los datos observados da un valor de su verosimilitud, el cual podemos optimizar variando $\boldsymbol{\theta}$ y σ^2 .

La suposición de independencia de los ruidos nos permite factorizar esta densidad en N términos separados, uno para cada dato:

$$L = p(\mathbf{y} | \mathbf{x}, \boldsymbol{\theta}, \sigma^2) = \prod_{i=1}^N p(y_i | \mathbf{x}_i, \boldsymbol{\theta}, \sigma^2) = \prod_{i=1}^N \mathcal{N}(\boldsymbol{\theta}^\top \mathbf{x}_i, \sigma^2)$$

Esto no quiere decir que los valores y_i son ellos mismos completamente independientes. Este no es el caso - los valores y_i , en promedio, varían con el \mathbf{x}_i , sugiriendo una clara dependencia entre ellos.

De hecho, son condicionalmente independientes - dado un valor para $\boldsymbol{\theta}$ (la parte determinista del modelo), los y_i son independientes.

Esta dependencia está encapsulada en el parámetro $\boldsymbol{\theta}$. Si conocemos $\boldsymbol{\theta}$, todo lo que queda son los errores entre los datos observados y $\boldsymbol{\theta}^\top \mathbf{x}_i$. Estos errores se asumen independientes. Por lo tanto, condicionado a $\boldsymbol{\theta}$, las observaciones son independientes.

Ahora mostraremos cómo podemos encontrar los valores de $\boldsymbol{\theta}$ y σ^2 que maximizan la verosimilitud.

4. Máxima verosimilitud

La verosimilitud nos da un único valor que nos indica cuán probable es nuestro conjunto de datos, dado el modelo actual (por modelo, nos referimos a la elección de $\boldsymbol{\theta}$ y σ^2). Como nuestro conjunto de datos es fijo, variar el modelo resultará en diferentes valores de verosimilitud. Una elección sensata de modelo sería aquella que maximice la verosimilitud. En otras palabras, seleccionaremos los parámetros del modelo que hagan nuestras observaciones más probables.

Sustituyendo la expresión para la función de densidad gaussiana:

$$\begin{aligned}\ln L &= \sum_{i=1}^N \ln \left(\frac{1}{\sqrt{2\pi\sigma^2}} \exp \left\{ -\frac{1}{2\sigma^2} (y_i - f(\mathbf{x}_i; \boldsymbol{\theta}))^2 \right\} \right) \\ &= \sum_{i=1}^N \left(-\frac{1}{2} \ln(2\pi) - \ln \sigma - \frac{1}{2\sigma^2} (y_i - f(\mathbf{x}_i; \boldsymbol{\theta}))^2 \right) \\ &= -\frac{N}{2} \ln 2\pi - N \ln \sigma - \frac{1}{2\sigma^2} \sum_{i=1}^N (y_i - f(\mathbf{x}_i; \boldsymbol{\theta}))^2\end{aligned}$$

Sustituyendo $f(\mathbf{x}_i; \boldsymbol{\theta}) = \boldsymbol{\theta}^\top \mathbf{x}_i$:

$$\ln L = -\frac{N}{2} \ln 2\pi - N \ln \sigma - \frac{1}{2\sigma^2} \sum_{i=1}^N (y_i - \boldsymbol{\theta}^\top \mathbf{x}_i)^2$$

5. Derivada respecto a $\boldsymbol{\theta}$

Podemos encontrar los parámetros óptimos tomando derivadas e igualándolas a cero. Para $\boldsymbol{\theta}$:

$$\frac{\partial \ln L}{\partial \boldsymbol{\theta}} = \frac{1}{\sigma^2} \frac{\partial}{\partial \boldsymbol{\theta}} \left[\sum_{i=1}^N (y_i - \boldsymbol{\theta}^\top \mathbf{x}_i)^2 \right] = \mathbf{0}$$

Notar que $\frac{\partial \ln L}{\partial \boldsymbol{\theta}}$ es un vector, por lo que lo igualamos a $\mathbf{0}$, un vector de ceros del mismo tamaño.

La solución de máxima verosimilitud para $\boldsymbol{\theta}$ coincide por lo tanto con la solución que ya hemos derivado para el caso de mínimos cuadrados. Minimizar la pérdida cuadrática es equivalente a la solución de máxima verosimilitud si se asume que el ruido es gaussiano.:

$$\hat{\boldsymbol{\theta}} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}$$

Recordar la notación matricial que utilizamos antes:

$$\mathbf{X} = \begin{bmatrix} \mathbf{x}_1^\top \\ \mathbf{x}_2^\top \\ \vdots \\ \mathbf{x}_N^\top \end{bmatrix} = \begin{bmatrix} 1 & x_1^{(1)} & \cdots & x_1^{(D)} \\ 1 & x_2^{(1)} & \cdots & x_2^{(D)} \\ \vdots & \vdots & \vdots & \vdots \\ 1 & x_N^{(1)} & \cdots & x_N^{(D)} \end{bmatrix}, \quad \mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{bmatrix}$$

6. Derivada respecto a σ

Para obtener una expresión para σ^2 (asumiendo que $\boldsymbol{\theta} = \hat{\boldsymbol{\theta}}$) podemos seguir el mismo procedimiento. Tomando derivadas parciales e igualando a cero obtenemos:

$$\frac{\partial L}{\partial \sigma} = -\frac{N}{\sigma} + \frac{1}{\sigma^3} \sum_{i=1}^N \left(y_i - \mathbf{x}_i^\top \hat{\boldsymbol{\theta}} \right)^2 = 0.$$

Reorganizando, obtenemos $\hat{\sigma}^2$, la estimación de máxima verosimilitud para σ^2 :

$$\hat{\sigma}^2 = \frac{1}{N} \sum_{i=1}^N \left(y_i - \mathbf{x}_i^\top \hat{\boldsymbol{\theta}} \right)^2$$

La varianza es simplemente el error cuadrático medio.

En notación matricial, usando el hecho de que

$$\sum_{i=1}^N \left(y_i - \mathbf{x}_i^\top \hat{\boldsymbol{\theta}} \right)^2 = (\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\theta}})^\top (\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\theta}})$$

tenemos:

$$\begin{aligned} \hat{\sigma}^2 &= \frac{1}{N} (\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\theta}})^\top (\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\theta}}) \\ &= \frac{1}{N} \left(\mathbf{y}^\top \mathbf{y} - 2\mathbf{y}^\top \mathbf{X}\hat{\boldsymbol{\theta}} + \hat{\boldsymbol{\theta}}^\top \mathbf{X}^\top \mathbf{X}\hat{\boldsymbol{\theta}} \right) \end{aligned}$$

Sustituyendo $\hat{\boldsymbol{\theta}} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}$ y observando que

$$\hat{\boldsymbol{\theta}}^\top = \mathbf{y}^\top \mathbf{X} (\mathbf{X}^\top \mathbf{X})^{-1}$$

pues $(\mathbf{X}^\top \mathbf{X})^{-1}$ es simétrica:

$$\begin{aligned} \hat{\sigma}^2 &= \frac{1}{N} \left(\mathbf{y}^\top \mathbf{y} - 2\mathbf{y}^\top \mathbf{X} (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y} + \mathbf{y}^\top \mathbf{X} (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{X} (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y} \right) \\ &= \frac{1}{N} \left(\mathbf{y}^\top \mathbf{y} - \mathbf{y}^\top \mathbf{X} (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y} \right) \\ &= \frac{1}{N} \left(\mathbf{y}^\top \mathbf{y} - \mathbf{y}^\top \mathbf{X}\hat{\boldsymbol{\theta}} \right) \end{aligned}$$

7. Modelos más complejos y sobreajuste

Sustituyendo la expresión para $\hat{\sigma}^2$ en la expresión de la log-verosimilitud nos da el valor de la log-verosimilitud en el máximo:

$$\begin{aligned}\ln L &= -\frac{N}{2} \ln 2\pi - \frac{N}{2} \ln \hat{\sigma}^2 - \frac{1}{2\hat{\sigma}^2} N \hat{\sigma}^2 \\ &= -\frac{N}{2} (1 + \ln 2\pi) - \frac{N}{2} \ln \hat{\sigma}^2.\end{aligned}$$

Esto nos indica que el valor máximo de L seguirá aumentando a medida que disminuimos $\hat{\sigma}^2$.

La varianza del ruido incorporado en el modelo para capturar efectos que la parte determinista de nuestro modelo (es decir, $f(\mathbf{x}; \boldsymbol{\theta})$) no puede. Una forma de disminuir σ^2 es modificar $f(\mathbf{x}; \boldsymbol{\theta})$ de manera que pueda capturar más de la variabilidad en los datos, es decir, hacerlo más flexible. El riesgo es caer en sobreajuste.

8. Incertidumbre en las estimaciones de los parámetros

El valor que obtenemos para $\hat{\boldsymbol{\theta}}$ está fuertemente influenciado por los valores particulares del ruido en los datos. Queremos saber cuánta incertidumbre hay en $\hat{\boldsymbol{\theta}}$.

Recordar lo que $\boldsymbol{\theta}$ y $\hat{\boldsymbol{\theta}}$ significan. Hemos hipotetizado un modelo responsable de generar los datos:

$$y_i = \boldsymbol{\theta}^\top \mathbf{x}_i + \epsilon_i$$

donde $\boldsymbol{\theta}$ representa el verdadero valor de los parámetros y ϵ_i es una variable aleatoria con distribución normal. Entonces la distribución generadora de targets

$$p(\mathbf{y} | \mathbf{X}, \boldsymbol{\theta}, \sigma^2)$$

es un producto de densidades normales:

$$p(\mathbf{y} | \mathbf{X}, \boldsymbol{\theta}, \sigma^2) = \prod_{i=1}^N p(y_i | \mathbf{x}_i, \boldsymbol{\theta}, \sigma^2) = \prod_{i=1}^N \mathcal{N}(\boldsymbol{\theta}^\top \mathbf{x}_i, \sigma^2)$$

Un producto de densidades Gaussianas univariadas puede escribirse como una densidad Gaussiana multivariada con una covarianza diagonal:

$$p(\mathbf{y} | \mathbf{X}, \boldsymbol{\theta}, \sigma^2) = \mathcal{N}(\mathbf{X}\boldsymbol{\theta}, \sigma^2 \mathbf{I}).$$

Ahora, $\hat{\boldsymbol{\theta}}$ es una estimación del verdadero valor del parámetro $\boldsymbol{\theta}$. Calculando la esperanza de $\hat{\boldsymbol{\theta}}$ con respecto a la distribución generadora nos dirá qué esperamos que sea $\hat{\boldsymbol{\theta}}$, en promedio:

$$\begin{aligned}\mathbb{E}_{\mathbf{y}|\mathbf{X}}\{\hat{\boldsymbol{\theta}}\} &= (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbb{E}_{\mathbf{y}|\mathbf{X}}\{\mathbf{y}\} \\ &= (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{X} \boldsymbol{\theta} \\ &= \boldsymbol{\theta}\end{aligned}$$

Este resultado nos indica que el valor esperado de nuestra aproximación $\hat{\boldsymbol{\theta}}$ es el verdadero valor del parámetro. Esto significa que nuestro estimador no tiene sesgo.

La variabilidad en la estimación de $\hat{\boldsymbol{\theta}}$ se obtiene de su matriz de covarianza. Esta matriz de covarianza nos proporciona dos piezas de información útiles:

- Las entradas diagonales son las varianzas de los coeficientes individuales de $\hat{\boldsymbol{\theta}}$, nos indican cuánta variabilidad podríamos esperar en los parámetros individuales, es decir, qué tan bien están definidos por los datos.
- Las entradas fuera de la diagonal nos indican cómo covarían los parámetros. Si los valores son altos y positivos, nos dice que aumentar uno requerirá un aumento en el otro para mantener un buen modelo. Valores negativos grandes nos dicen lo contrario: aumentar uno causará una disminución en el otro. Valores cercanos a cero nos indican que los parámetros no dependen uno del otro.

La matriz de covarianza es:

$$\begin{aligned}\text{cov}\{\hat{\boldsymbol{\theta}}\} &= \mathbb{E}_{\mathbf{y}|\mathbf{X}}\{\hat{\boldsymbol{\theta}}\hat{\boldsymbol{\theta}}^\top\} - \mathbb{E}_{\mathbf{y}|\mathbf{X}}\{\hat{\boldsymbol{\theta}}\}\mathbb{E}_{\mathbf{y}|\mathbf{X}}\{\hat{\boldsymbol{\theta}}\}^\top \\ &= \mathbb{E}_{\mathbf{y}|\mathbf{X}}\{\hat{\boldsymbol{\theta}}\hat{\boldsymbol{\theta}}^\top\} - \boldsymbol{\theta}\boldsymbol{\theta}^\top\end{aligned}$$

donde hemos usado la esperanza de $\hat{\boldsymbol{\theta}}$ que derivamos anteriormente. El primer término es

$$\begin{aligned}\mathbb{E}_{\mathbf{y}|\mathbf{X}}\{\hat{\boldsymbol{\theta}}\hat{\boldsymbol{\theta}}^\top\} &= \mathbb{E}_{\mathbf{y}|\mathbf{X}}\left\{\left((\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}\right)\left((\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}\right)^\top\right\} \\ &= (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbb{E}_{\mathbf{y}|\mathbf{X}}\{\mathbf{y}\mathbf{y}^\top\} \mathbf{X} (\mathbf{X}^\top \mathbf{X})^{-1}\end{aligned}$$

La covarianza de \mathbf{y} es, por definición, $\sigma^2 \mathbf{I}$ y su media es $\mathbf{X}\boldsymbol{\theta}$. Luego tenemos:

$$\text{cov}\{\mathbf{y}\} = \sigma^2 \mathbf{I} = \mathbb{E}_{\mathbf{y}|\mathbf{X}}\{\mathbf{y}\mathbf{y}^\top\} - \mathbb{E}_{\mathbf{y}|\mathbf{X}}\{\mathbf{y}\}\mathbb{E}_{\mathbf{y}|\mathbf{X}}\{\mathbf{y}\}^\top$$

Por lo tanto, podemos reorganizar esta expresión para obtener:

$$\begin{aligned}\mathbb{E}_{\mathbf{y}|\mathbf{X}}\{\mathbf{y}\mathbf{y}^\top\} &= \mathbb{E}_{\mathbf{y}|\mathbf{X}}\{\mathbf{y}\}\mathbb{E}_{\mathbf{y}|\mathbf{X}}\{\mathbf{y}\}^\top + \sigma^2\mathbf{I} \\ &= \mathbf{X}\boldsymbol{\theta}(\mathbf{X}\boldsymbol{\theta})^\top + \sigma^2\mathbf{I} \\ &= \mathbf{X}\boldsymbol{\theta}\boldsymbol{\theta}^\top\mathbf{X}^\top + \sigma^2\mathbf{I}\end{aligned}$$

Finalmente obtenemos:

$$\mathbb{E}_{\mathbf{y}|\mathbf{X}}\{\boldsymbol{\theta}\boldsymbol{\theta}^\top\} = \boldsymbol{\theta}\boldsymbol{\theta}^\top + \sigma^2(\mathbf{X}^\top\mathbf{X})^{-1}$$

Obtenemos así la expresión para la covarianza de $\widehat{\boldsymbol{\theta}}$:

$$\begin{aligned}\text{cov}\{\widehat{\boldsymbol{\theta}}\} &= \boldsymbol{\theta}\boldsymbol{\theta}^\top + \sigma^2(\mathbf{X}^\top\mathbf{X})^{-1} - \boldsymbol{\theta}\boldsymbol{\theta}^\top \\ &= \sigma^2(\mathbf{X}^\top\mathbf{X})^{-1}\end{aligned}$$

Para $\widehat{\sigma}^2$ los cálculos son un poco más engorrosos, pero el resultado es:

$$\mathbb{E}_{\mathbf{y}|\mathbf{X}}\{\widehat{\sigma}^2\} = \sigma^2\left(1 - \frac{D}{N}\right)$$

donde D es el número de atributos (el número de columnas de \mathbf{X}).

Asumiendo que $D < N$, es decir, el número de atributos que medimos para cada observación es menor que el número de datos, entonces nuestra estimación de la varianza será, en promedio, menor que la varianza real:

$$\mathbb{E}_{\mathbf{y}|\mathbf{X}}\{\widehat{\sigma}^2\} < \sigma^2.$$

A diferencia de $\widehat{\boldsymbol{\theta}}$, este estimador es sesgado.

9. Variabilidad en las predicciones

Además de obtener una indicación de la variabilidad de nuestra estimación del parámetro, $\widehat{\boldsymbol{\theta}}$, tiene interés proporcionar indicaciones de cualquier variabilidad o incertidumbre en nuestras predicciones.

Supongamos que observamos un nuevo conjunto de atributos, $\mathbf{x}_{\text{nuevo}}$. Nos gustaría predecir la salida y_{nuevo} y además, la variabilidad asociada con esta salida, σ_{nuevo}^2 .

Para predecir y_{nuevo} , multiplicamos $\mathbf{x}_{\text{nuevo}}$ por $\widehat{\boldsymbol{\theta}}$: $y_{\text{nuevo}} = \widehat{\boldsymbol{\theta}}^\top \mathbf{x}_{\text{nuevo}}$.

Para empezar podemos calcular su esperanza:

$$\begin{aligned}\mathbb{E}_{\mathbf{y}|\mathbf{X}} \{y_{\text{nuevo}}\} &= \mathbb{E}_{\mathbf{y}|\mathbf{X}} \{\widehat{\boldsymbol{\theta}}\}^\top \mathbf{x}_{\text{nuevo}} \\ &= \boldsymbol{\theta}^\top \mathbf{x}_{\text{nuevo}}.\end{aligned}$$

El valor esperado de nuestra predicción es el nuevo atributo de datos multiplicado por el verdadero $\boldsymbol{\theta}$.

La varianza es:

$$\sigma_{\text{nuevo}}^2 = \text{var} \{y_{\text{nuevo}}\} = \mathbb{E}_{\mathbf{y}|\mathbf{X}} \{y_{\text{nuevo}}^2\} - (\mathbb{E}_{\mathbf{y}|\mathbf{X}} \{y_{\text{nuevo}}\})^2$$

Al sustituir $y_{\text{nuevo}} = \widehat{\boldsymbol{\theta}}^\top \mathbf{x}_{\text{nuevo}}$:

$$\begin{aligned}\text{var} \{y_{\text{nuevo}}\} &= \mathbb{E}_{\mathbf{y}|\mathbf{X}} \left\{ \left(\widehat{\boldsymbol{\theta}}^\top \mathbf{x}_{\text{nuevo}} \right)^2 \right\} - (\boldsymbol{\theta}^\top \mathbf{x}_{\text{nuevo}})^2 \\ &= \mathbb{E}_{\mathbf{y}|\mathbf{X}} \left\{ \mathbf{x}_{\text{nuevo}}^\top \widehat{\boldsymbol{\theta}} \widehat{\boldsymbol{\theta}}^\top \mathbf{x}_{\text{nuevo}} \right\} - \mathbf{x}_{\text{nuevo}}^\top \boldsymbol{\theta} \boldsymbol{\theta}^\top \mathbf{x}_{\text{nuevo}}.\end{aligned}$$

Sustituyendo la expresión para $\widehat{\boldsymbol{\theta}}$:

$$\text{var} \{y_{\text{nuevo}}\} = \sigma^2 \mathbf{x}_{\text{nuevo}}^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{x}_{\text{nuevo}}$$

Notar que al sustituir nuestra expresión para $\text{cov}\{\widehat{\boldsymbol{\theta}}\}$ la varianza puede reescribirse como

$$\sigma_{\text{nuevo}}^2 = \mathbf{x}_{\text{nuevo}}^\top \text{cov}\{\widehat{\boldsymbol{\theta}}\} \mathbf{x}_{\text{nuevo}}.$$

Para resumir, nuestra predicción y la varianza asociada:

$$\begin{aligned}y_{\text{nuevo}} &= \mathbf{x}_{\text{nuevo}}^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{x}^\top \mathbf{y} = \mathbf{x}_{\text{nuevo}}^\top \widehat{\boldsymbol{\theta}} \\ \sigma_{\text{nuevo}}^2 &= \sigma^2 \mathbf{x}_{\text{nuevo}}^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{x}_{\text{nuevo}}.\end{aligned}$$

σ^2 es la verdadera varianza del ruido del conjunto de datos. En su lugar, podemos usar nuestra estimación, $\widehat{\sigma}^2$.

10. Varianza de y_{nuevo} en el caso univariado

Recordar que en el caso univariado la matriz $\mathbf{X}^\top \mathbf{X}$ es:

$$\mathbf{X}^\top \mathbf{X} = \begin{bmatrix} 1 & 1 & \cdots & 1 \\ x_1 & x_2 & \cdots & x_N \end{bmatrix} \begin{bmatrix} 1 & x_1 \\ 1 & x_2 \\ \vdots & \vdots \\ 1 & x_N \end{bmatrix} = \begin{bmatrix} N & \sum_{i=1}^N x_i \\ \sum_{i=1}^N x_i & \sum_{i=1}^N x_i^2 \end{bmatrix}$$

y la inversa de la matriz es:

$$\begin{aligned}
 (\mathbf{X}^\top \mathbf{X})^{-1} &= \frac{1}{N \sum_{i=1}^N x_i^2 - \left(\sum_{i=1}^N x_i\right)^2} \begin{bmatrix} \sum_{i=1}^N x_i^2 & -\sum_{i=1}^N x_i \\ -\sum_{i=1}^N x_i & N \end{bmatrix} \\
 &= \frac{1}{\overline{x^2} - \bar{x}^2} \begin{bmatrix} \overline{x^2} & -\bar{x} \\ -\bar{x} & 1 \end{bmatrix}
 \end{aligned}$$

Luego, la varianza de la predicción:

$$\begin{aligned}
 \sigma_{\text{nuevo}}^2 &= \frac{\sigma^2}{s_x^2} [1, x_{\text{nuevo}}] \begin{bmatrix} s_x^2 + \bar{x}^2 & -\bar{x} \\ -\bar{x} & 1 \end{bmatrix} \begin{bmatrix} 1 \\ x_{\text{nuevo}} \end{bmatrix} \\
 &= \frac{\sigma^2}{s_x^2} [1, x_{\text{nuevo}}] \begin{bmatrix} s_x^2 + \bar{x}^2 - \bar{x}x_{\text{nuevo}} \\ x_{\text{nuevo}} - \bar{x} \end{bmatrix} \\
 &= \frac{\sigma^2}{s_x^2} (s_x^2 + (x_{\text{nuevo}} - \bar{x})^2) \\
 &= \sigma^2 \left(1 + \frac{(x_{\text{nuevo}} - \bar{x})^2}{s_x^2} \right)
 \end{aligned}$$