

Modelo lineal Bayesiano

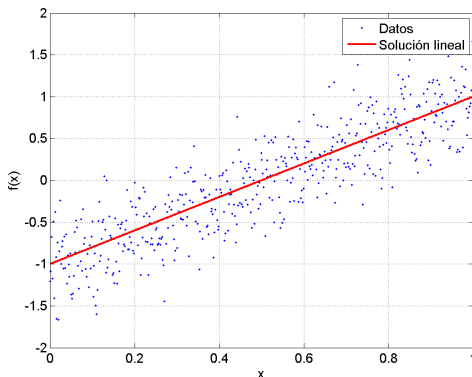
A. M. Alvarez-Meza, Ph.D.
amalvarezme@unal.edu.co

Departamento de ingeniería eléctrica, electrónica y computación
Universidad Nacional de Colombia-sede Manizales



- 1 Modelo lineal y mínimos cuadrados
- 2 Modelos Bayesianos
- 3 Máxima verosimilitud
- 4 Modelo lineal Bayesiano
- 5 Maximum a posteriori (MAP)

Modelo lineal



- $y = f(x) = mx + b \Rightarrow$ noción lineal desde algebra básica.
- $\mathbf{y} = f(\mathbf{x}) = \mathbf{xw} + \mathbf{b} \Rightarrow$ extensión algebra vectorial.
- Los datos no siempre llegan limpios y no siempre comparten relaciones lineales!

Modelo lineal: extensión matricial

$$\hat{\mathbf{y}} = f(\mathbf{X}) = \mathbf{X}\mathbf{w} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{bmatrix} = \begin{bmatrix} x_{11} & x_{12} & \dots & x_{1P} \\ x_{21} & x_{22} & \dots & x_{2P} \\ \vdots & \vdots & \ddots & \vdots \\ x_{N1} & x_{N2} & \dots & x_{NP} \end{bmatrix} \begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_P \end{bmatrix}$$

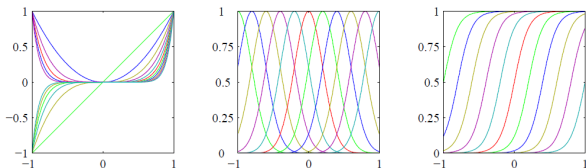
- $\hat{\mathbf{y}} \in \mathbb{R}^N$, $\mathbf{X} \in \mathbb{R}^{N \times P}$, $\mathbf{w} \in \mathbb{R}^P$
- N : # muestras.
- P : # características.

Cómo encontrar los parámetros del modelo lineal (\mathbf{w})?

$$\epsilon(\mathbf{w}, \lambda) = \frac{1}{2} \|\mathbf{e}\|_2^2 + \frac{1}{2} \lambda \|\mathbf{w}\|_2^2$$

- $\|\mathbf{e}\|_2^2$: cuantifica el desajuste entre los datos y el modelo lineal de aproximación.
- $\|\mathbf{w}\|_2^2$: cuantifica el sobreajuste (complejidad) de la solución.
- $\lambda \in \mathbb{R}^+$: parámetro de balance ("trade-off").
- Se necesita un λ que garantice una solución simple pero "exacta".

Regresión desde representación no lineal



- Polinomial: $\phi(\mathbf{x}) = [\mathbf{x}^j]_{j=1}^D$, D : grado del polinomio.
- Exponencial: $\phi(\mathbf{x}|\{\boldsymbol{\mu}_j\}_{j=1}^Q, \sigma) = \left[\exp\left(\frac{-\|\mathbf{x} - \boldsymbol{\mu}_j\|_2^2}{2\sigma^2}\right) \right]_{j=1}^Q$.
- Sigmoidal:
 $\phi(\mathbf{x}|\{\boldsymbol{\mu}_j\}_{j=1}^Q, \sigma) = \left[1 / (1 + \exp(\|\mathbf{x} - \boldsymbol{\mu}_j\|_2^2 / (2\sigma^2))) \right]_{j=1}^Q$.

Solución del modelo lineal sobre representaciones no lineales

- Solución sobre \mathbf{X} :

$$\mathbf{w}^* = \left(\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I} \right)^{-1} \mathbf{X}^\top \mathbf{y}$$

$$y_i = \mathbf{x}_i \mathbf{w}$$

- Solución sobre Φ :

$$\mathbf{w}^* = \left(\Phi^\top \Phi + \lambda \mathbf{I} \right)^{-1} \Phi^\top \mathbf{y}$$

$$\mathbf{X} \in \mathbb{R}^{N \times P}, \mathbf{w} \in \mathbb{R}^P, \mathbf{X}^\top \mathbf{X} \in \mathbb{R}^{P \times P}.$$

$$y_i = \phi(\mathbf{x}_i) \mathbf{w}$$

$$\Phi \in \mathbb{R}^{N \times Q}, \mathbf{w} \in \mathbb{R}^Q, \Phi^\top \Phi \in \mathbb{R}^{Q \times Q}.$$

Thomas Bayes



Thomas Bayes

1701–1761

Thomas Bayes was born in Tunbridge Wells and was a clergyman as well as an amateur scientist and a mathematician. He studied logic and theology at Edinburgh University and was elected Fellow of the Royal Society in 1742. During the 18th century, issues regarding probability arose in connection with

gambling and with the new concept of insurance. One particularly important problem concerned so-called inverse probability. A solution was proposed by Thomas Bayes in his paper 'Essay towards solving a problem in the doctrine of chances', which was published in 1764, some three years after his death, in the *Philosophical Transactions of the Royal Society*. In fact, Bayes only formulated his theory for the case of a uniform prior, and it was Pierre-Simon Laplace who independently rediscovered the theory in general form and who demonstrated its broad applicability.

Máxima verosimilitud I

- Supongamos y dado como

$$y = f(\mathbf{x}, \mathbf{w}) + \epsilon$$

donde $\epsilon \sim \mathcal{N}(0, \beta^{-1})$.

- La incertidumbre en y está dada como

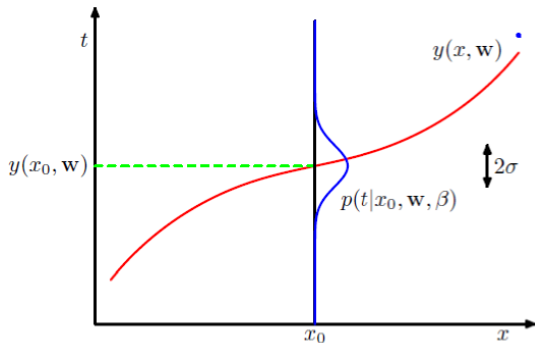
$$p(y|\mathbf{x}, \mathbf{w}, \beta^{-1}) = \mathcal{N}(y|f(\mathbf{x}, \mathbf{w}), \beta^{-1}).$$

- Consideremos un conjunto de datos (de entrenamiento)

$$\mathbf{X} = \begin{bmatrix} \mathbf{x}_1^\top, \dots, \mathbf{x}_N^\top \end{bmatrix}^\top, \\ \mathbf{y} = [y_1, \dots, y_N]^\top,$$

con $\mathbf{x}_i \in \mathbb{R}^P$, $y \in \mathbb{R}$.

Ejemplo:



Máxima verosimilitud II

- Suponiendo que los datos son independientes e idénticamente distribuidos (iid) y $f(\mathbf{x}, \mathbf{w}) = \phi(\mathbf{x})\mathbf{w}$:

$$p(\mathbf{y}|\mathbf{X}, \mathbf{w}, \beta) = \prod_{i=1}^N \mathcal{N}(y_i | \phi(\mathbf{x}_i)\mathbf{w}, \beta^{-1})$$

- Tomando el logaritmo de la verosimilitud se tiene

$$\begin{aligned} \log(p(\mathbf{y}|\mathbf{X}, \mathbf{w}, \beta)) &= \sum_{i=1}^N \log(\mathcal{N}(y_i | \phi(\mathbf{x}_i)\mathbf{w}, \beta^{-1})) \\ &= \frac{N}{2} \log(\beta) - \frac{N}{2} \log(2\pi) - \beta E_D(\mathbf{w}) \end{aligned}$$

$$E_D(\mathbf{w}) = \frac{1}{2} \sum_{i=1}^N (y_i - \phi(\mathbf{x}_i)\mathbf{w})^2$$

- Cuáles son los \mathbf{w} y el parámetro β que mejor explican los datos?

Máxima verosimilitud III

- Maximizar la verosimilitud es equivalente a minimizar $-\beta E_D(\mathbf{w})$
- De nuevo,

$$\begin{aligned} E_D(\mathbf{w}) &= \frac{1}{2} \sum_{i=1}^N (y_i - \phi(\mathbf{x}_i)\mathbf{w})^2 \\ &= \frac{1}{2} (\mathbf{y} - \Phi\mathbf{w})^\top (\mathbf{y} - \Phi\mathbf{w}) \end{aligned}$$

donde

$$\Phi = \begin{bmatrix} \phi(\mathbf{x}_1) \\ \phi(\mathbf{x}_2) \\ \vdots \\ \phi(\mathbf{x}_N) \end{bmatrix}$$

$$\Phi \in \mathbb{R}^{N \times Q}, \phi : \mathbb{R}^P \rightarrow \mathbb{R}^Q.$$

Máxima verosimilitud IV

- La verosimilitud logarítmica está dada entonces como:

$$\log(p(\mathbf{y}|\mathbf{X}, \mathbf{w}, \beta)) = \frac{N}{2} \log(\beta) - \frac{N}{2} \log(2\pi) - \frac{\beta}{2} (\mathbf{y} - \Phi \mathbf{w})^\top (\mathbf{y} - \Phi \mathbf{w})$$

- Se tiene entonces

$$\begin{aligned} \frac{\partial \log(p(\mathbf{y}|\mathbf{X}, \mathbf{w}, \beta))}{\partial \mathbf{w}} &= -\frac{\beta}{2} \frac{\partial}{\partial \mathbf{w}} [(\mathbf{y} - \Phi \mathbf{w})^\top (\mathbf{y} - \Phi \mathbf{w})] \\ &= -\frac{\beta}{2} \frac{\partial}{\partial \mathbf{w}} [\mathbf{y}^\top \mathbf{y} - 2\mathbf{y}^\top \Phi \mathbf{w} + \mathbf{w}^\top \Phi^\top \Phi \mathbf{w}] \end{aligned}$$

- Note que

$$\frac{\partial (\mathbf{a}^\top \mathbf{x})}{\partial \mathbf{x}} = \mathbf{a}, \quad \frac{\partial (\mathbf{x}^\top \mathbf{A} \mathbf{x})}{\partial \mathbf{x}} = (\mathbf{A} + \mathbf{A}^\top) \mathbf{x}.$$

Máxima verosimilitud V

- Por ende

$$\frac{\partial \log(p(\mathbf{y}|\mathbf{X}, \mathbf{w}, \beta))}{\partial \mathbf{w}} = \beta \left[\Phi^\top \mathbf{y} \mathbf{w} - \Phi^\top \Phi \mathbf{w} \right].$$

- La solución de máxima verosimilitud para \mathbf{w} está dada como:

$$\mathbf{w}_{ML} = \left(\Phi^\top \Phi \right)^{-1} \Phi^\top \mathbf{y},$$

- La solución de máxima verosimilitud para β se obtiene de:

$$\frac{\partial \log(p(\mathbf{y}|\mathbf{X}, \mathbf{w}, \beta))}{\partial \beta} = \frac{N}{2\beta} - \frac{1}{2} (\mathbf{y} - \Phi \mathbf{w})^\top (\mathbf{y} - \Phi \mathbf{w}).$$

- Y así,

$$\frac{1}{\beta_{ML}} = \frac{1}{N} (\mathbf{y} - \Phi \mathbf{w}_{ML})^\top (\mathbf{y} - \Phi \mathbf{w}_{ML}) = \frac{1}{N} \sum_{i=1}^N (y_i - \phi(\mathbf{x}_i))^2.$$

- Al igual que en mínimos cuadrados, el modelo de máxima verosimilitud puede regularizarse.
- Se pretende controlar el sobre entrenamiento.
- La función de error toma la forma

$$E_D(\mathbf{w}) + \lambda E_W(\mathbf{w}) = \frac{1}{2} \sum_{i=1}^N (y_i - \phi(\mathbf{x}_i)\mathbf{w})^2 + \frac{\lambda}{2} \mathbf{w}^\top \mathbf{w},$$

donde $E_W(\mathbf{w}) = \frac{1}{2} \mathbf{w}^\top \mathbf{w}$.

- El valor de \mathbf{w} que minimiza $E_D(\mathbf{w}) + \lambda E_W(\mathbf{w})$ está dado por:

$$\mathbf{w} = \left(\Phi^\top \Phi + \lambda I \right)^{-1} \Phi^\top \mathbf{y}.$$

Modelo Bayesiano como alternativa a la regularización

- Modelos lineales basados en mínimos cuadrados y máxima verosimilitud pueden regularizarse.
- El valor del parámetro de regularización se impone o asume.
- Una alternativa a la regularización es el tratamiento Bayesiano.

- La verosimilitud del modelo lineal esta dado por:

$$p(\mathbf{y}|\phi, \mathbf{w}, \beta) = \mathcal{N}(\mathbf{y}|\phi\mathbf{w}, \beta^{-1}\mathbf{I}).$$

- En máxima verosimilitud se calculó una solución puntual para $\mathbf{w} \rightarrow \mathbf{w}_{ML}$.
- En estimación Bayesiana se asume un prior para \mathbf{w} y se calcula la probabilidad a posteriori de \mathbf{w} dados los datos \mathbf{y} .
- El posteriori sobre \mathbf{w} se usa para hacer predicciones.

Teorema de Bayes

- Para calcular el posterior sobre \mathbf{w} se usa el teorema de Bayes:

$$p(\mathbf{w}|\mathbf{y}) = \frac{p(\mathbf{y}|\mathbf{w}) p(\mathbf{w})}{p(\mathbf{y})},$$

donde $p(\mathbf{y})$ es la evidencia, $p(\mathbf{y}|\mathbf{w})$ es la verosimilitud y $p(\mathbf{w})$ es el prior.

- Usando el modelo $\mathbf{y} = f(\mathbf{w}, \mathbf{x}) + \epsilon$; con $\epsilon \sim \mathcal{N}(0, \beta^{-1})$, la verosimilitud es conocida.
- Dependiendo del prior que se escoja para \mathbf{w} , es posible calcular analíticamente el posterior.

Prior y posterior

- Asumiendo que el prior es Gaussiano, el posterior es igualmente Gaussiano.
- En particular, suponga que $p(\mathbf{w}) = \mathcal{N}(\mathbf{w}|\mathbf{m}_0, \mathbf{S}_0)$.
- Usando propiedades de la Gaussiana, se puede demostrar que:

$$\begin{aligned} p(\mathbf{w}|\mathbf{y}) &= \frac{p(\mathbf{y}|\mathbf{w}) p(\mathbf{w})}{p(\mathbf{y})} = \frac{\mathcal{N}(\mathbf{y}|\Phi\mathbf{w}, \beta^{-1}\mathbf{I}) \mathcal{N}(\mathbf{w}|\mathbf{m}_0, \mathbf{S}_0)}{p(\mathbf{y})} \\ &= \mathcal{N}(\mathbf{w}|\mathbf{m}_N, \mathbf{S}_N). \end{aligned}$$

donde

$$\begin{aligned} \mathbf{m}_N &= \mathbf{S}_N \left(\mathbf{S}_0^{-1} \mathbf{m}_0 + \beta \Phi^\top \mathbf{y} \right) \\ \mathbf{S}_N^{-1} &= \mathbf{S}_0^{-1} + \beta \Phi^\top \Phi. \end{aligned}$$

Tarea: Demuestre las igualdades anteriores.

Propiedades de la Gaussiana

- Dadas una distribución Gaussiana marginal para \mathbf{x} , y una distribución Gaussiana condicional para \mathbf{y} , dado \mathbf{x} , de la forma:

$$p(\mathbf{x}) = \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Delta}^{-1})$$
$$p(\mathbf{y}|\mathbf{x}) = \mathcal{N}(\mathbf{y}|\mathbf{A}\mathbf{x} + \mathbf{b}, \mathbf{L}^{-1})$$

la distribución marginal de \mathbf{y} , y la distribución condicional de \mathbf{x} dado \mathbf{y} están dadas por:

$$p(\mathbf{y}) = \mathcal{N}(\mathbf{y}|\mathbf{A}\boldsymbol{\mu} + \mathbf{b}, \mathbf{L}^{-1} + \mathbf{A}\boldsymbol{\Delta}^{-1}\mathbf{A}^\top)$$
$$p(\mathbf{x}|\mathbf{y}) = \mathcal{N}(\mathbf{x}|\boldsymbol{\Sigma}\{\mathbf{A}^\top\mathbf{L}(\mathbf{y} - \mathbf{b}) + \boldsymbol{\Delta}\boldsymbol{\mu}\}, \boldsymbol{\Sigma})$$

donde

$$\boldsymbol{\Sigma} = \left(\boldsymbol{\Delta} + \mathbf{A}^\top\mathbf{L}\mathbf{A}\right)^{-1}$$

Prior más simple

- Un prior más sencillo sigue la forma $p(\mathbf{w}) = \mathcal{N}(\mathbf{w}|\mathbf{0}, \alpha^{-1}\mathbf{I})$.
- El posterior está dado como

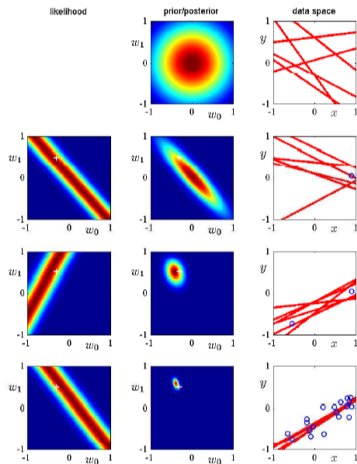
$$p(\mathbf{w}|\mathbf{y}) = \mathcal{N}(\mathbf{w}|\mathbf{m}_N, \mathbf{S}_N).$$

donde

$$\begin{aligned}\mathbf{m}_N &= \beta \mathbf{S}_N \Phi^\top \mathbf{y} \\ \mathbf{S}_N^{-1} &= \alpha \mathbf{I} + \beta \Phi^\top \Phi.\end{aligned}$$

Tarea: Demuestre las igualdades anteriores.

Ejemplo: posterior



$$\beta^{-1} = 0.04, \alpha = 2, w_0 = -0.3, w_1 = 0.5.$$

Maximum a posteriori (MAP)

- La regularización se puede ver como la estimación del Maximum A Posteriori (MAP).
- El logaritmo del posterior es una función de \mathbf{w}

$$\log(p(\mathbf{w}|\mathbf{y})) = -\frac{\beta}{2} \sum_{i=1}^N (y_i - \phi(\mathbf{x}_i)\mathbf{w})^2 - \frac{\alpha}{2} \mathbf{w}^\top \mathbf{w} + \text{cte}$$

- Equivalente a la regularización si $\lambda = \alpha/\beta$.

Tarea: Demuestre la igualdad anterior.

Distribución predictiva

- **Objetivo:** hacer predicciones de y para nuevos valores \mathbf{x} .
- Denotemos ese nuevo valor de entrada como \mathbf{x}_* y la predicción resultante como y_* .
- La distribución predictiva para y_* está dada como:

$$p(y_*|\mathbf{y}, \alpha, \beta, \mathbf{x}_*) = \int p(y_*|\mathbf{w}, \beta, \mathbf{x}_*) p(\mathbf{w}|\mathbf{y}, \alpha, \beta) d\mathbf{w}$$

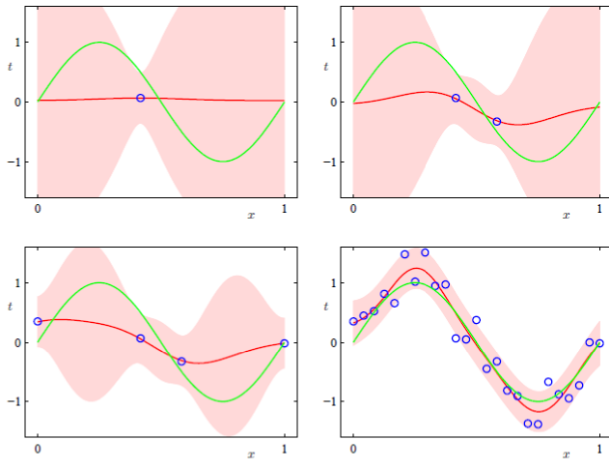
- Usando las propiedades de la Gaussiana se puede demostrar que:

$$p(y_*|\mathbf{y}, \alpha, \beta, \mathbf{x}_*) = \mathcal{N}(y_*|\phi(\mathbf{x}_*)\mathbf{m}_N, \sigma_N^2(\mathbf{x}_*)) ,$$

$$\text{donde } \sigma_N^2(\mathbf{x}_*) = \beta^{-1} + \phi(\mathbf{x}_*)\mathbf{S}_N\phi(\mathbf{x}_*)^\top$$

- β y α se han asumido conocidos.

Ejemplo: distribución predictiva



Aproximación de la evidencia I

- Si no se conocen α y β , cómo se pueden estimar a partir del conjunto de entrenamiento?
- En un tratamiento Bayesiano general, se ponen priors sobre α y β y se calculan los posteriores.
- Alternativamente, se puede estimar como los parámetros que maximizan la evidencia $p(\mathbf{y}|\alpha, \beta)$.
- Este método se conoce como máxima verosimilitud tipo II, aproximación de la evidencia, Bayes empírico.

Aproximación de la evidencia II

- La evidencia está dada como

$$p(\mathbf{y}|\alpha, \beta) = \int p(\mathbf{y}|\mathbf{w}, \beta) p(\mathbf{w}|\alpha) d\mathbf{w}$$

$$\text{evidencia} = \int \text{verosimilitud} \times \text{prior}$$

- Maximizando la expresión anterior en función de α y β :

$$\alpha = \frac{\gamma}{\mathbf{m}_N^\top \mathbf{m}_N},$$

donde $\gamma = \sum_i \frac{\lambda_i}{\alpha + \lambda_i}$, $(\beta \Phi^\top \Phi) \mathbf{u}_i = \lambda_i \mathbf{u}_i$,

$$\mathbf{m}_N = \beta (\beta \Phi^\top \Phi + \alpha \mathbf{I})^{-1} \Phi^\top \mathbf{y}.$$

- Note que es una solución iterativa para α , porque γ y \mathbf{m}_N dependen de α .

- Derivando $\log(p(\mathbf{y}|\alpha, \beta))$ con respecto a β

$$\frac{1}{\beta} = \frac{1}{N - \gamma} |\mathbf{y} - \Phi \mathbf{m}_N|^2.$$

- De nuevo esta es una solución implícita para β , porque \mathbf{m}_N depende de β , la solución es iterativa.

En un cuaderno (notebook) responda a las siguientes preguntas con ejemplos concretos de implementación sobre Python. Envíe/comparta su notebook al correo amalvarezme@unal.edu.co.

- Consultar el funcionamiento (modelo matemático, función de costo y optimización) de los algoritmos de regresión Bayesiana `sklearn.linear_model.BayesianRidge` y `sklearn.linear_model.ARDRegression` según sus implementaciones en el paquete Scikit-Learn de Python.
- Consulte cómo podría incluir funciones de representación no lineal en los algoritmos de regresión del punto anterior.
- Utilizando las bases de datos estudiadas en la sesión de mínimos cuadrados, realice un análisis comparativo de los métodos de regresión mediante validación cruzada. Recuerde sintonizar los parámetros libres de cada algoritmo.



Hansen, P. C. (1998).

Rank-deficient and discrete ill-posed problems: numerical aspects of linear inversion.
(Vol. 4). Siam.



Bishop, C. (2006).

Pattern recognition.
Machine Learning, 128.