

# Apuntes Gaussian Processes

Patricio Whittingslow

March 2019

## 1 Prefacio

Es importante la palabra lineal. Aquí no hay funciones con cuadrados ni nada por el estilo. Existe una función global lineal que toma como argumento las variables elegidas.

### Glosario Términos

$\mathbf{x}$  Input vector ( $\mathbb{D}$ ). Que variables eligo para modelo.

$\mathbf{y}$  Vector objetivo ( $n$ ). Observaciones/mediciones.

$\mathbf{w}$  Weight vector ( $\mathbb{D}$ ). Aquí van los parámetros a obtener de la regresión lineal. La “solución” de la regresión lineal.

$\mathbb{D}$  Dimensión del problema. Cuantas variables elegí para el modelo.

$n$  Cantidad de observaciones.

$X$  Matriz de diseño ( $\mathbb{D} \times n$ ).

$\varepsilon \sim \mathcal{N}(0, \sigma_n^2)$  Ruido. Sigue una distribución gaussiana independiente, idénticamente distribuida con promedio 0 y varianza  $\sigma_n^2$

$|\mathbf{v}| = (\sum_i \mathbf{v}_i^2)^{\frac{1}{2}}$  Longitud euclidiana del vector  $\mathbf{v}$ .

$\sigma_n$  Asumimos que las mediciones están alejadas de la función  $f$  por causa de ruido  $\varepsilon$ .

$\Sigma_p$  Matriz de covarianza o **matriz de varianza-covarianza**.

### Glosario Subíndices

$i, j, k, p, q$  Refiere a un elemento de un tensor, sea vector, matriz, etc.

## 1.1 Introducción

Dados datos de entrenamiento  $(X, \mathbf{y})$  se quiere efectuar predicciones para nuevas entradas  $\mathbf{x}_*$  usando una función  $f$ , por lo tanto, es claro que el problema es *inductivo*. Para lograr esto tenemos que suponer ciertas características sobre nuestra función  $f$ . Si no acotáramos  $f$  entonces cualquier función que sea consistente con los datos de entrenamiento valdría. Este es el problema que intenta resolver el aprendizaje de máquina, el **Machine Learning**.

## 2 Regresión Lineal

Para que sea más didáctico la aplicación de la regresión lineal se va a plantear un problema a resolver: *Se tiene un cilindro de diámetro  $D$  y largo  $L$  que enfrenta un fluido de viscosidad  $\mu$ , densidad  $\rho$  y velocidad  $U_\infty$ . Obtener una expresión para la fuerza  $F_D$  que el fluido ejerce sobre el cilindro?*

Cabe destacar que el problema no es lineal y que probablemente la regresión lineal que obtengamos sea válida para unos pocos puntos. Igual intentaremos hallar la regresión lineal.

Primero tenemos que partir de datos, preferiblemente **muchos** datos. Vamos a obviar la variable  $L$ , ya que si dividimos la fuerza por el largo obtenemos la fuerza por unidad de largo y reducimos la dimensión del problema sin perder información.

Abajo está la matriz de diseño para el problema. Se efectuaron  $n$  observaciones o mediciones, por ende también se tienen  $n$  resultados para la fuerza  $y = F_{D_i}$ .

$$X = \begin{bmatrix} \mu_1 & \mu_2 & \dots & \mu_n \\ \rho_1 & \rho_2 & \dots & \rho_n \\ U_1 & U_2 & \dots & U_n \\ D_1 & D_2 & \dots & D_n \end{bmatrix}$$
$$\mathbf{y} = \begin{Bmatrix} F_1 \\ F_2 \\ \vdots \\ F_n \end{Bmatrix}$$

### 2.1 No-Bayesiana

Como las mediciones nunca son exactas (siempre hay error humano o por el sistema de medición) se puede proponer ruido ( $\varepsilon$ ) de distribución gaussiana  $\varepsilon \sim \mathcal{N}(0, \sigma_n^2)$  que se suma a nuestra regresión para obtener los valores de  $y$  tal que

$$f(\mathbf{x}) = \mathbf{x}^T \mathbf{w} \quad y = f(\mathbf{x}) + \varepsilon$$

Se presentaron dos nuevos vectores al usuario,  $\mathbf{x}$  y  $\mathbf{w}$ , los vectores *input* y de *pesos* (*weight* en inglés), respectivamente. Ambos tienen dimensión del

problema, en este caso  $\mathbb{D} = 4$ .

$$\mathbf{x} = \begin{Bmatrix} \mu_i \\ \rho_i \\ U_i \\ D_i \end{Bmatrix} \quad \mathbf{w} = \begin{Bmatrix} w_a \\ w_b \\ w_c \\ w_d \end{Bmatrix}$$

Behold. La regresión lineal entonces quedaría de la forma

$$f(\mathbf{x}) = w_a \mu_i + w_b \rho_i + w_c U_i + w_d D_i = y_i = F_i$$

Resolver el problema (efectivamente: obtener el vector  $\mathbf{w}$ ) requiere un uso pesado de la probabilidad, en particular el teorema de Bayes. También se tiene que elegir un *prior*, es decir, dejar expresadas nuestras creencias de los parámetros (pesos) antes de mirar las observaciones. Esto toma la forma de la matriz de covarianza  $\Sigma_p$ . Si suponemos que los parámetros se distribuyen según  $\mathbf{w} \sim \mathcal{N}(\mathbf{0}, \Sigma_p)$  entonces:<sup>1</sup>

$$p(\mathbf{w}|X, \mathbf{y}) \sim \mathcal{N}(\bar{\mathbf{w}}, A)$$

$$A = (\sigma_n^{-2} X X^T + \Sigma_p^{-1})^{-1} \quad (1)$$

$$\bar{\mathbf{w}} = \sigma_n^{-2} (\sigma_n^{-2} X X^T + \Sigma_p^{-1})^{-1} X \mathbf{y} \quad (2)$$

Donde  $\bar{\mathbf{w}}$  son los promedios de los pesos hallados.  $A$  es la *matriz de covarianza de la distribución posterior*.

Se suele referir al promedio de los pesos  $\bar{\mathbf{w}}$  como MAP, *maximum a posteriori*.

## 2.2 Bayesiana

Se ha hablado de la formulación No-Bayesiana hasta ahora. De que se trata esto de si es Bayesiano o no? En la formulación Bayesiana el caso prueba tiene la forma del vector de entrada, y la regresión se calcula para ese vector prueba. Por ende, en el esquema Bayesiano no existe un único vector  $\mathbf{w}$ . Un subíndice  $*$  se refiere a un caso prueba.

$$p(f_*|\mathbf{x}_*, X, \mathbf{y}) = \mathcal{N}(\sigma_n^{-2} \mathbf{x}_*^T A X \mathbf{y}, \mathbf{x}_*^T A \mathbf{x}_*) \quad (3)$$

Se usará esta formulación para la regresión en el espacio funcional.  $A$  se definió anteriormente.

## 3 Regresión en el espacio funcional

Si bien hemos obtenido una regresión cuando calculamos  $\bar{\mathbf{w}}$ , no resuelve el problema debido a que el problema *no es lineal*.

---

<sup>1</sup>Aún no se ha explicado como se calcula  $\Sigma_p$ .

Esto se resuelve proyectando las variables de entrada a un espacio funcional. Nuestra función para la regresión tendría la forma

$$f(\mathbf{x}) = \phi(\mathbf{x})^T \mathbf{w} \quad (4)$$

esta función  $\phi$  mapea el vector de entrada  $\mathbf{x}$  de un espacio  $\mathbb{D}$  a un espacio funcional  $N$ . La aplicación de este modelo es análoga a la regresión lineal excepto que donde antes teníamos  $X$  ahora tenemos:

$$\Phi = \Phi(X)$$

$$p(f_* | \mathbf{x}_*, X, \mathbf{y}) \sim \mathcal{N}(\sigma_n^{-2} \phi(\mathbf{x}_*)^T A \Phi \mathbf{y}, \phi(\mathbf{x}_*)^T A \phi(\mathbf{x}_*)) \quad (5)$$

donde

$$A = (\sigma_n^{-2} \Phi \Phi^T + \Sigma_p^{-1})^{-1}$$

### 3.1 Kernel

Si uno quisiese ahorrar tiempo y no calcular la inversa de  $A$  (calculo computacional pesado) podría recurrir a la expresión (donde  $\phi(x) = \phi$ ):

$$p(f_* | \mathbf{x}_*, X, \mathbf{y}) \sim \mathcal{N}(\phi_*^T \Sigma_p \Phi (K + \sigma_n^2 I)^{-1} \mathbf{y}, \phi_*^T \Sigma_p \phi_* - \phi_*^T \Sigma_p \Phi (K + \sigma_n^2 I)^{-1} \Phi^T \Sigma_p \phi_*) \quad (6)$$

donde  $K = \Phi^T \Sigma_p \Phi$ .

Note que en la ecuación 6 es prevalente el producto de el espacio funcional en las formas  $\phi_*^T \Sigma_p \Phi$ ,  $\phi_*^T \Sigma_p \phi_*$  y  $\Phi^T \Sigma_p \Phi$ . Definiendo entonces una función  $\psi(\mathbf{x}) = \Sigma_p^{-\frac{1}{2}} \phi(\mathbf{x})$  obtenemos

$$k(\mathbf{x}, \mathbf{x}') = \psi(\mathbf{x}) \cdot \psi(\mathbf{x}') \quad (7)$$

donde  $\mathbf{x}$  y  $\mathbf{x}'$  pueden ser las variables de entrada o casos prueba (denotado con subíndice  $*$ ). Se le suele decir la función kernel<sup>2</sup> a  $k$  y es de sumo interés en la rama del aprendizaje de maquina. Si reescribiéramos 6:

$$p \sim \mathcal{N}(k(\mathbf{x}_*, X)(k(X, X) + \sigma_n^2 I)^{-1} \mathbf{y}, k(\mathbf{x}_*, \mathbf{x}_*) - k(\mathbf{x}_*, X)(k(X, X) + \sigma_n^2 I)^{-1} k(X, \mathbf{x}_*)) \quad (8)$$

Puede ser que llegado a este punto el lector se sienta ansioso por no tener las herramientas de como aplicar 8. Para aliviar dichas dudas se presenta una posible función de covarianza, la *squared exponential* (SE).<sup>3</sup>

$$\text{cov}(f(\mathbf{x}_p), f(\mathbf{x}_q)) = k(\mathbf{x}_p, \mathbf{x}_q) = \exp(-\frac{1}{2} |\mathbf{x}_p - \mathbf{x}_q|^2 / \ell) \quad (9)$$

donde  $\ell$  es la longitud característica del proceso gaussiano.

<sup>2</sup>También llamada función covarianza.

<sup>3</sup>también llamada la gaussiana o *Radial Basis Function*

**Definición 3.1** *Un Proceso Gaussiano es una colección de cualquier número finito de variables aleatorias, todas distribuidas gaussianamente.*

### 3.2 Predicción con ruido

Suponiendo que  $y = f(x) + \varepsilon$  siendo el ruido  $\varepsilon \sim \mathcal{N}(0, \sigma_n^2)$

$$\text{cov}(y_p, y_q) = k(\mathbf{x}_p, \mathbf{x}_q) + \sigma_n^2 \delta_{pq} \quad o \quad \text{cov}(\mathbf{y}) = K(X, X) + \sigma_n^2 I$$

La distribución de  $\mathbf{y}$  y las salidas según el *prior* elegido son

$$\begin{bmatrix} \mathbf{y} \\ \mathbf{f} \end{bmatrix} \sim \mathcal{N} \left( \mathbf{0}, \begin{bmatrix} K(X, X) + \sigma_n^2 I & K(X, X_*) \\ K(X_*, X) & K(X_*, X_*) \end{bmatrix} \right)$$