

# CURSO SOBRE MODELAMIENTO ESTADÍSTICO Y SISTEMAS RECOMENDADORES

## Clase Presencial 2: Modelos Estadísticos y Big DATA

## 2.1.- MODELOS ESTADÍSTICOS

# Modelos estadísticos

- Los modelos estadísticos se construyen en base a modelos de probabilidad.
- Los datos se tratan como realizaciones de variables aleatorias  $\mathbf{X}_1, \dots, \mathbf{X}_n$ , donde  $\mathbf{X}_i$  es en si mismo un vector de variables aleatorias recogidas para la  $i$ -ésima unidad experimental en una muestra de tamaño  $n$  desde una población de interés.
- El supuesto es que  $\mathbf{X} = (\mathbf{X}_1, \dots, \mathbf{X}_n)$  son generados desde un modelo de probabilidad  $P$ .

# Modelo estadístico

- Si  $P$  se conoce completamente, por ejemplo, al conocer la densidad  $f(\mathbf{x})$ , entonces no existe la necesidad de usar estadística.
- El “problema” estadístico surge cuando existe incertidumbre sobre  $P$ .
- Los modelos estadísticos surgen cuando  $P$  se trata como un miembro de una familia  $\mathcal{M} = \{P_\theta : \theta \in \Theta\}$  indexada por un conjunto de parámetros  $\theta$ , que a su vez pertenece a un conjunto  $\Theta$  (espacio paramétrico).

## 2.2.- TIPOS DE MODELOS ESTADÍSTICOS

# Modelos paramétricos

- Modelos que se describen a través de un número finito de, típicamente, valores reales se denominan como **modelos paramétricos o finito dimensionales**.
- Los modelos paramétricos pueden ser descritos por la familia

$$\mathcal{M} = \{P_{\theta} : \theta \in \Theta \subset \mathbb{R}^p\}$$

donde la dimensión  $p$  es un entero positivo.

# Modelos paramétricos

- En el caso de modelos paramétricos, el problema se especifica de tal forma que el valor de los parámetros (o alguna función de ellos) es de importancia para los investigadores.
- El objetivo del análisis estadístico es entonces especificar un valor posible de  $\theta$ , o determinar un subconjunto de  $\Theta$  para el cuál es posible afirmar que contiene o no a  $\theta$ , en base a la muestra observada.

# Modelos paramétricos

Ejemplo 2:

Suponga que  $X_1, \dots, X_n \stackrel{i.i.d.}{\sim} N(\mu, \sigma^2)$ .



# Modelos paramétricos

Ejemplo 3 (regresión lineal):

Suponga que  $Y_i \mid x_i \stackrel{\text{ind.}}{\sim} N(\beta_0 + \beta_1 x_i, \sigma^2)$ ,  $i = 1, \dots, n$ .

# Modelos no-paramétricos

- En muchas situaciones restringir las inferencias a formas paramétricas específicas pueden limitar los tipos de conclusiones que se pueden obtener.
- Nos gustaría relajar supuestos paramétricos para ganar mayor flexibilidad y robustez en contra de la especificación incorrecta de un modelo paramétrico.
- En estos casos, podemos considerar modelos donde la clase de distribuciones de probabilidad es tan grande que el parámetro  $\theta$  es de dimensión infinita.

# Modelos no-paramétricos

## Ejemplo 4 (regresión):

Suponga que  $Y_i | x_i \stackrel{\text{ind.}}{\sim} N(m(x_i), \sigma^2)$ ,  $i = 1, \dots, n$ , donde  $m : \mathbb{R} \rightarrow \mathbb{R}$  es una función continua.

# Modelos no-paramétricos

## Ejemplo 5 (regresión):

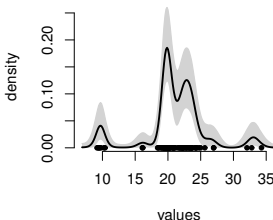
Suponga que  $Y_i | x_i \stackrel{\text{ind.}}{\sim} f(\cdot - m(x_i))$ ,  $i = 1, \dots, n$ , donde  $m : \mathbb{R} \rightarrow \mathbb{R}$  es una función continua y  $f$  es una densidad continua en todas partes y definida en  $\mathbb{R}$ .

## 2.3.- BIG DATA

- El tamaño muestral y el número de variables es enorme.
- Los datos no caben en la memoria o no pueden ser almacenados en una sola máquina.



- Los datos deben ser super raros también.



# Ejemplos

- Las distribuciones de datos de transacciones en Internet tienen picos en cero y en otros valores discretos (por ejemplo, 1 o \$ 99).
- Colas grandes que importan (por ejemplo, \$ 12 mil/mes de gasto de usuario de eBay).
- El espacio de características potenciales es inmanejablemente grande.
- No podemos escribir modelos simples para explicar los datos.

- Siempre tenemos que almacenar todos los datos para aprender sobre  $\theta$ ?
- Podemos resumir los datos y todavía tener la misma información sobre  $\theta$ ?
- Consideremos un ejemplo simple:





# Estadísticos suficientes

## Definición:

Un estadístico  $T(\mathbf{X})$  es **suficiente para  $\theta$**  si la distribución condicional de la muestra  $\mathbf{X}$  dado el valor de  $T(\mathbf{X})$  no depende de  $\theta$ .

# Estadísticos suficientes

Ejemplo:

Sean  $X_1, \dots, X_n \stackrel{i.i.d.}{\sim} \text{Bernoulli}(\theta)$ ,  $\theta \in (0, 1)$ .  $T(\mathbf{X}) = \sum_{i=1}^n X_i$  es un estadístico suficiente para  $\theta$ .

# Estadísticos suficientes

## Teorema de factorización:

Sea  $f(\mathbf{x} \mid \theta)$  la función de densidad o de probabilidad de la muestra  $\mathbf{X}$ . Un estadístico  $T(\mathbf{X})$  es suficiente para  $\theta$  si y sólo si existen unas funciones  $g(t \mid \theta)$  y  $h(\mathbf{x})$ , tal que para todo punto muestral  $\mathbf{x}$  y todo punto en el espacio paramétrico  $\theta$ ,

$$f(\mathbf{x} \mid \theta) = g(T(\mathbf{x}) \mid \theta) h(\mathbf{x}).$$

# Estadísticos suficientes

Ejemplo:

Sean  $X_1, \dots, X_n \stackrel{i.i.d.}{\sim} N(\mu, \sigma^2)$ , donde  $\theta = (\mu, \sigma^2)$ .  $T(\mathbf{X}) = (\bar{X}, S^2)$  es un estadístico suficiente para  $\theta$ .

# Estadísticos mínimo suficientes

## Definición:

Un estadístico suficiente  $T(\mathbf{X})$  para  $\theta$  se denomina **mínimo suficiente**, si para cualquier otro estadístico suficiente  $T'(\mathbf{X})$ ,  $T(\mathbf{X})$  es una función de  $T'(\mathbf{X})$ .

# Comentarios finales

- Una reducción de los datos no es siempre posible y útil. Existe alguna esperanza si  $\Theta \subseteq \mathbb{R}^p$ , con  $p$  finito.
- No existe necesidad de Big DATA en esos casos!!!
- Una reducción es inútil cuando  $\Theta$  es un espacio de dimensión infinita (i.e., modelos noparamétricos).

## Comentarios finales

- La promesa del Big Data no tiene que ver con conocer algunas características de una población de interés con mayor precisión.
- El interés en los grandes conjuntos de datos se debe a que ellos pueden ser “extraños” y nos permiten conocer los mecanismos complejos que los generan (modelos noparamétricos).
- En estos contextos, puede ser difícil o incluso contraproducente emplear modelos estadísticos simples para el proceso de aprendizaje, donde el número de parámetros se fija a priori (modelos paramétricos).

# Comentarios finales

- Los modelos Bayesianos noparamétricos permiten esta flexibilidad (los modelos pueden crecer en tamaño y complejidad con la llegada de más datos), y por un tratamiento coherente de las incertidumbres.

