

Estadística (1937)  
**Unidad 7.**  
**Introducción a la ciencia de datos**

Marcelo Ruiz

Departamento de Matemática, Universidad Nacional de Río Cuarto

2023

¿Qué es la “ciencia de datos”?  
Aprendizaje estadístico  
La maldición de la dimensionalidad

## Aprendizaje Estadístico

# Dato

Todo “dato”:

- es una construcción.
- surge como respuesta a preguntas en el contexto de un problema <sup>1</sup>.



---

<sup>1</sup> Imagen tomada de <https://twitter.com/Yyomepregunto2/photo>

## El problema determina el dato

Afirma Chalmers<sup>2</sup> que:

se consiguen hechos relevantes midiendo la concentración de ozono en varios lugares, mientras que no se logra nada midiendo la longitud de los cabellos de 105 jóvenes de Sidney.



---

<sup>2</sup>Chalmers, A. (1982). *¿Qué es esa cosa llamada ciencia?. Siglo XXI.*

## Provisoriamente

¿De qué hablamos cuando hablamos de ciencia de datos?

Hastie, en 2015, sostiene: <sup>3</sup>

- El **Aprendizaje Automático** (Machine Learning) construye algoritmos que pueden aprender de los datos.
- El **Aprendizaje Estadístico** (Statistical Learning) es una rama de la estadística aplicada que emergió en respuesta al AA, enfatizando en los modelos estadísticos y en la evaluación de la incertidumbre.
- La **Ciencia de Datos** (CD) es la extracción de conocimiento de los datos, utilizando ideas de matemática, estadística, AA, ingeniería,...

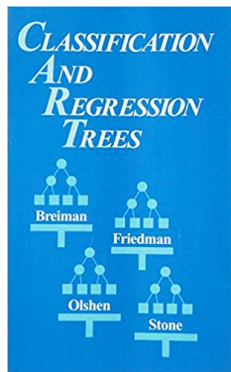
---

<sup>3</sup> [https://web.stanford.edu/~hastie/TALKS/SLBD\\_new.pdf](https://web.stanford.edu/~hastie/TALKS/SLBD_new.pdf)

# Las transformaciones científico-técnica y nuevas dinámicas mundiales

Me interesa mencionar algunos aspectos de su emergencia, mucho tiempo atrás.

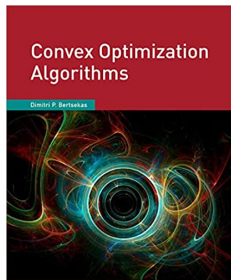
- Ciertas mutaciones en dinámicas disciplinares clásicas y novedosas emergencias:
  - en líneas de la computación y de la estadística.
  - en lingüística (de Peirce a IBM).
  - en la ingeniería de Pattern Recognition.
- Problemas
  - de la industria
  - en inteligencia y defensa (DARPA).



## La dinámica de la financiarización y la conectividad

Más recientemente:

- En términos disciplinares, las múltiples articulaciones entre
  - estadística
  - teoría de aproximación.
  - programación convexa
  - ecuaciones diferenciales
  - álgebra lineal numérica
  - optimización
  - teoría de probabilidad
  - ingeniería de software, etc.
- Las fuerzas fácticas de las dinámicas financieras y la conectividad.



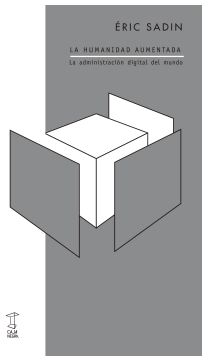
Ed. del 2015

## La administración digital del mundo

Hay que inscribir la enumeración anterior en:

- una administración robotizada de la existencia humana,
- una duplicación digital del mundo,
- una administración digital del mundo y,
- en una mutación decisiva en nuestro vínculo con la técnica

en la perspectiva que plantea Eric Sadin (2018), cuyo texto fue traducido por Javier Blanco.





## Las culturas generativa y predictiva

Si los datos vienen generados por:



en relación al [análisis de datos](#), Breiman<sup>4</sup> distingue dos objetivos:

- **Predecir** cuál será el valor de la respuesta  $Y$  para valores futuros de  $\mathbf{X}$ .
- **Inferir** cómo la naturaleza está asociando la variable respuesta  $Y$  a las variables de entrada  $\mathbf{X}$ .

<sup>4</sup>Breiman, L. (2001), "Statistical Modeling: the Two Cultures", Statistical Science, 16, 199–231.

## Las culturas generativa y predictiva

- En la **cultura generativa**:
  - se asume que hay un verdadero modelo generando los datos y,
  - el análisis de datos requiere de la inferencia.
- En la **cultura predictiva**<sup>5</sup>
  - el modelo dice poco del mecanismo subyacente que genera los datos y,
  - se hace foco en la precisión de la predicción (machine learning).

---

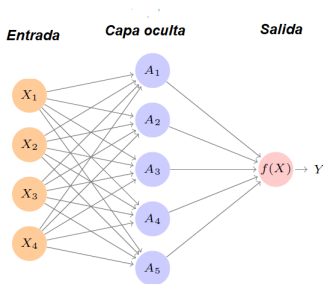
<sup>5</sup> Maldonado, J., Picco, M. Ruiz, M. (2022). Comparación de métodos de clasificación desde la perspectiva de la predicción y la selección del modelo. Ponencia en el Congreso de la SAE

## Las culturas generativa y predictiva

- En la cultura generativa:

$$Y = \sum_{i=1}^p \beta_i X_i + \epsilon, \quad \epsilon \sim N(\mu, \sigma^2).$$

- En la cultura predictiva:



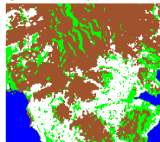
## Predicción con Random Forest

**Aprendizaje supervisado.** Los sensores de imágenes hiperespectrales proveen cientos de anchos de banda del espectro electromagnético de la misma área sobre la superficie terrestre.

- A cada pixel le corresponde  $\mathbf{X} = (X_1, \dots, X_p)$ ; cada entrada  $X_i$  expresa la reflectancia en una longitud de onda específica.
- Los niveles de la variable respuesta: agua, cemento y dos tipos de vegetación.
- El objetivo: entrenar un clasificador de tal modo que, a partir de un nuevo dato,  $x_{\text{nuevo}}$ , **prediga (clasifique)** bien.



Imagen real. Hyperion EO-1.

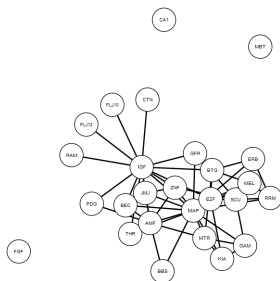


Mapa temático.  
Random forest.

## Inferencia con Modelos gráficos

**Aprendizaje no supervisado.** Ciertos tipos de cáncer quedan caracterizados por perfiles génicos:

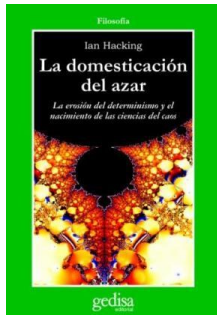
- Cada entrada,  $X_i$ , del vector  $\mathbf{X} = (X_1, \dots, X_p)$  representa la expresión de la actividad del gen  $i$ ,  $i = 1, \dots, p$ .
- Interesa estimar la asociación de dos genes condicional a los restantes genes, para todos los pares de genes posibles.
- Si dos genes están relacionados condicionalmente a los restantes ponemos un lado el grafo.
- El objetivo: **inferir** el grafo asociado a  $\mathbf{X}$ .



Grafo estimado para 26 genes asociados al cáncer de mama de pacientes con respuesta patológica completa.

## Alta dimensionalidad

- ¿Los “Big data” caracterizan el surgimiento de la CD?
- Los datos censales surgieron siglos atrás y son big data.
- En el análisis de datos clásico  $n$  es grande en relación al número de variables  $p$ .



¡Muy buen texto de historia!

## Alta dimensionalidad

- En las últimas décadas, la **novedad** práctica y teórica es

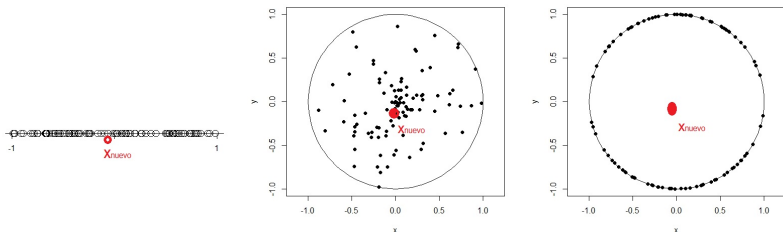
$$p \gg n \text{ ó } p \approx n$$
$$\mathbb{X} = \begin{pmatrix} & X_1 & X_2 & \cdots & X_p \\ \text{caso 1} & x_{11} & x_{12} & \cdots & x_{1p} \\ \text{caso 2} & x_{21} & x_{22} & \cdots & x_{2p} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ \text{caso } n & x_{n1} & x_{n2} & \cdots & x_{np} \end{pmatrix}$$

y  $p \rightarrow \infty$ .

- En genómica, el número de genes  $p = 20000$  y  $n = 80$ .
- En las redes sociales, el número de usuarios puede ser mucho menor que el número de variables que escoge.
- En teledetección, en imágenes hiperespectrales  $p \approx n$ .

## Alta dimensionalidad

En alta dimensionalidad los los puntos están aislados



Entornos para  $p = 1, 2$  y  $p$  “grande”.



## Datos administrativos y transaccionales

David Hand<sup>6</sup> puso de relevancia aspectos importantes de la recolección y análisis de los datos administrativos (DA).

- Los DA se recolectan por sistemas administrativos incluyendo el sector público.
- Son generados durante el curso de alguna operación y retenidos en una base de datos.
- Los DA han sido principales impulsores del desarrollo de la tecnología de minería de datos y de la CD.

---

<sup>6</sup>Hand, D. (2018). Statistical challenges of administrative and transaction data. J. R. Statist. Soc. A 181, Part 3, pp. 555–605

## Datos administrativos y transaccionales

Ejemplos de DA<sup>7</sup>:

- Registros educativos
- Registros hospitalarios tomados de la visita de pacientes.
- Datos de los clientes de entidades financieras y bancarias.
- Información recolectada en forma automática de las prestaciones en obras sociales.

---

<sup>7</sup> Sitio oficial de la Unión Europea:

[https://ec.europa.eu/eurostat/cros/content/administrative-data-0\\_en](https://ec.europa.eu/eurostat/cros/content/administrative-data-0_en)

## Ejemplos de DA

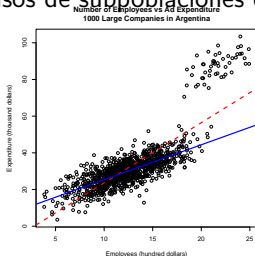
Algunas características:

- No son muestras estadísticas, suele ser relevada la población completa.
- Se necesitan medidas descriptivas (funcionales).
- Los DA necesitan ser
  - resumidos
  - descriptos
  - analizadospara descubrir diferentes niveles de complejidad posiblemente ocultos en diferentes niveles de complejidad y dimensionalidad.
- En general no hay problemas con la cantidad de observaciones sino con la calidad.

## Análisis de datos robusto

Los DA pueden estar ingresados incorrectamente, conteniendo:<sup>8</sup>.

- valores atípicos por celda; p. ej. 11.1 en lugar de 1.11
- filas completas atípicas; p. ej., una muestra de sangre contaminada o la inclusión de casos de subpoblaciones diferentes.



Small example to illustrate performance of robust regression algorithms.

<sup>8</sup> Ruiz, M., Yohai, V., Zamar, R. (2018). Discussion of: “Statistical challenges of administrative and

## Algunas tareas que se realizan en Ciencia de Datos

Donoho (2015) enumera grupos de tareas:

- Recolección, preparación y exploración.
  - Recolección: desde el clásico diseño experimental a las técnicas actuales de manipulación<sup>9</sup>
  - Preparación, como curación de datos, entre otras tareas.
  - Exploración de datos, en el sentido clásico.
- Representación y transformación de los datos.
  - Bases de datos: implica conocer estructuras, transformaciones, algoritmos involucrados.
  - Representaciones matemáticas: como transformada de Fourier, wavelets, multiescala (en deep learning).
- Cómputos. Incluye aportar nuevas bibliotecas a lenguajes de comunidades (como R o Python), utilización de clusters.

---

<sup>9</sup> web scraping, Pubmed scraping, 38 image processing, and Twitter, Facebook, and Reddit mining

## ¿Qué tareas se realizan en Ciencia de Datos?

- Visualización y presentación. Va desde histogramas, boxplots, bagplots hasta dispositivos más complejos como cellhandlers, etc.
- Modelización de los datos.
  - Modelización generativa: es la estadística matemática académica tradicional.
  - Modelización predictiva: casi coincidente con el AA (machine learning).

## Áreas y subáreas de Estadística

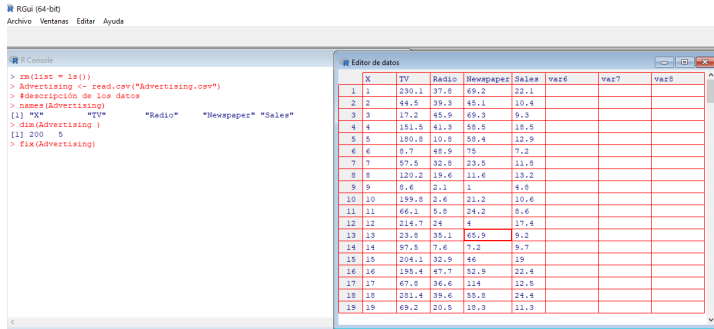
- Metodología estadística general.
  - Estimación
  - Inferencia
  - Selección de modelos
  - Análisis multivariado.
  - Métodos no paramétricos.
- Bioestadística y estadística del ambiente.
- Estadística en las ciencias físicas y en la industria.
- Estadística en economía y ciencias sociales.
- Estadística computacional y CD.
  - Algorítmica
  - Grandes datos (Big data)
  - Clasificación
  - Clustering
  - Análisis de datos.
  - Minería de datos
  - Visualización
  - Análisis de imágenes
  - Aprendizaje automático
  - Datos de alta dimensión
  - Aprendizaje estadístico
  - Lenguajes y software
  - Pattern recognition

¿Qué es la “ciencia de datos”?  
Aprendizaje estadístico  
La maldición de la dimensionalidad

Primeras nociones  
¿Por qué estimar  $f$ ?  
¿Cómo estimar  $f$ ?  
Aprendizaje supervisado versus no supervisado  
Regresión versus Clasificación  
Precisión del modelo  
El compromiso sesgo-varianza  
El contexto de clasificación

## Problema motivador: estudio de mercado

### Datos “Advertising”



**Figura 1:**  $n = 200$  observaciones correspondientes, cada una, a un producto vendido en un mercado. Para cada producto conocemos el valor de venta y la inversión publicitaria en tres diferentes medios: tv, radio y diario impreso (James et al., 2021)



¿Qué es la “ciencia de datos”?  
Aprendizaje estadístico  
La maldición de la dimensionalidad

Primeras nociones

¿Por qué estimar  $f$ ?

¿Cómo estimar  $f$ ?

Aprendizaje supervisado versus no supervisado

Regresión versus Clasificación

Precisión del modelo

El compromiso sesgo-varianza

El contexto de clasificación

## Estudio de mercado

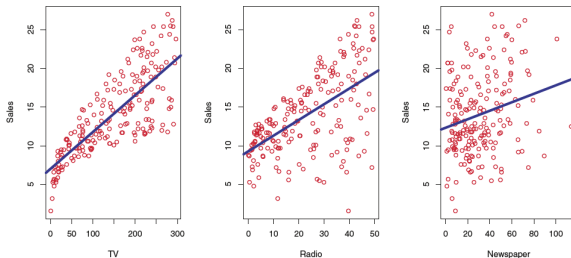


Figura 2: Datos “Advertising”. Representación de ventas (**sales**) versus inversión según medio publicitario (**TV**, **Radio** y **Newspaper**)

## Estudio de mercado

- No se pueden incrementar el valor de venta del producto pero sí se puede controlar el gasto en publicidad.
- Si determinamos que existe una asociación entre publicidad y ventas entonces podríamos recomendarle al vendedor cómo ajustar los montos publicitarios.

### Objetivo

Desarrollar un modelo que pueda ser utilizado para predecir ventas sobre la base de tres presupuestos en medios

¿Qué es la “ciencia de datos”?  
Aprendizaje estadístico  
La maldición de la dimensionalidad

Primeras nociones

¿Por qué estimar  $f$ ?

¿Cómo estimar  $f$ ?

Aprendizaje supervisado versus no supervisado

Regresión versus Clasificación

Precisión del modelo

El compromiso sesgo-varianza

El contexto de clasificación

## Estudio de mercado

### Variables del modelo

- *Variables de entrada (inputs), características (features) o predictores: “presupuesto” en TV, Radio y Newspaper.*
- *Variable output, de salida, respuesta o variable dependiente: “valor de venta” o sales.*

## Modelo de regresión

Más formalmente, contamos con una variable aleatoria cuantitativa  $Y$ , denominada **respuesta**, y un vector de **predictoras**<sup>10</sup>  $\mathbf{x} = (X_1, \dots, X_p)$  tal que

$$Y = f(\mathbf{x}) + \epsilon \quad (1)$$

donde  $\epsilon$ , una **variable aleatoria**, representa el **error**, del cual asumimos que posee media cero y es independiente de  $\mathbf{x}$ .

En esta formulación  $f$  representa la **información sistemática** que  $\mathbf{x}$  posee acerca de  $Y$ .

---

<sup>10</sup>Puede ser aleatorio o controlado

¿Qué es la “ciencia de datos”?  
**Aprendizaje estadístico**  
La maldición de la dimensionalidad

Primeras nociones

¿Por qué estimar  $f$ ?

¿Cómo estimar  $f$ ?

Aprendizaje supervisado versus no supervisado

Regresión versus Clasificación

Precisión del modelo

El compromiso sesgo-varianza

El contexto de clasificación

## ¿Por qué estimar $f$ ?

### ¿Por qué estimar $f$ ?

Dos razones:

- predicción
- inferencia

# Predicción

## Predicción

Se cuenta con un conjunto de inputs  $x$  pero el output no puede ser obtenido fácilmente.

## Ejemplo

- $X_1, \dots, X_p$  representan características de la sangre (muestra) de un paciente e  $Y$  es la respuesta, adversa, a un medicamento.
- $Y$  no se podrá obtener por experimentación si la adversidad es severa.  
Entonces queremos estimar  $Y$  utilizando las características.

## Predicción

$\hat{Y}$

- Como  $E(\epsilon) = 0$ , podemos predecir  $Y$  utilizando

$$\hat{Y} = \hat{f}(\mathbf{x})$$

con  $\hat{f}$  un estimador de  $f$ .

- En este contexto, no nos interesa la forma exacta de  $\hat{f}$  (caja negra) sino con qué exactitud predice  $Y$ .

## Predicción

La precisión depende de dos cantidades:

- *Error reducible*: asociado al tipo de estimador  $\hat{f}$ .
- *Error irreducible*: asociado al error  $\epsilon$  y no depende de la precisión de  $\hat{f}$ . Este error se debe a que puede contener variables no medidas.

Asumiendo  $\mathbf{x}$  y  $\hat{f}$  fijos:

$$E(Y - \hat{Y})^2 = \underbrace{|f(\mathbf{x}) - \hat{f}(\mathbf{x})|^2}_{\text{componente reducible}} + \underbrace{\text{VAR}(\epsilon)}_{\text{componente irreducible}} \quad (2)$$



¿Qué es la “ciencia de datos”?  
**Aprendizaje estadístico**  
La maldición de la dimensionalidad

Primeras nociones

¿Por qué estimar  $f$ ?

¿Cómo estimar  $f$ ?

Aprendizaje supervisado versus no supervisado

Regresión versus Clasificación

Precisión del modelo

El compromiso sesgo-varianza

El contexto de clasificación

# Inferencia

## Objetivo

Comprender la relación existente entre las variables input y output.

## Interrogantes

- ¿Qué predictores están asociados con qué respuestas?
- ¿Cuál es la relación entre la respuesta y cada predictor? Algunos predictores pueden provocar incrementos en  $Y$  y otros lo contrario, etc.
- ¿Puede la relación entre la respuesta y las predictoras resumirse adecuadamente en forma lineal o es más compleja?

## Ejemplos

### Predicción

Una compañía lanza una campaña de marketing directo.

- El objetivo es identificar las respuestas positivas a invitaciones por correo electrónico, basados en variables demográficas identificadas sobre cada individuo.
- No hay interés en comprender en profundidad la relación entre la respuesta a la campaña y las variables medidas sobre cada individuo.
- La compañía está interesada en obtener un modelo preciso para predecir la respuesta utilizando predictores.

## Ejemplos

### Inferencia

En el contexto de la base de datos “Advertising” el interés está centrado en responder a las siguientes preguntas:

- ¿Qué medio contribuye a las ventas?
- ¿Cuáles medios generan el mayor impulso de venta?
- ¿Cuál es el aumento en las ventas asociado a un incremento otorgado al presupuesto en TV?

## ¿Cómo estimar $f$ ?

Consideremos  $n$  realizaciones del modelo (1):

$$y_i = f(\mathbf{x}_i) + \epsilon_i, \quad i = 1, \dots, n \quad (3)$$

donde  $\mathbf{x}_i^T = (x_{i1}, \dots, x_{ip})$  es la  $i$ -ésima observación para la  $j$ -ésima predictora,  $i = 1, \dots, n$  y  $j = 1, \dots, p$ .

La matriz de datos es

$$\begin{pmatrix} y_1 & x_{11} & x_{12} & \dots & x_{1p} \\ y_2 & x_{21} & x_{22} & \dots & x_{2p} \\ \vdots & \vdots & & \vdots & \\ y_n & x_{n1} & x_{n2} & \dots & x_{np} \end{pmatrix} \quad (4)$$

¿Qué es la “ciencia de datos”?  
**Aprendizaje estadístico**  
La maldición de la dimensionalidad

Primeras nociones

¿Por qué estimar  $f$ ?

¿Cómo estimar  $f$ ?

Aprendizaje supervisado versus no supervisado

Regresión versus Clasificación

Precisión del modelo

El compromiso sesgo-varianza

El contexto de clasificación

## ¿Cómo estimar $f$ ?

### Métodos de estimación de $f$

- Métodos paramétricos.
- Métodos no-paramétricos.

## Métodos paramétricos

En estos métodos hay dos etapas

- a) Seleccionamos un modelo, por ejemplo lineal:

$$f(\mathbf{X}) = \sum_{i=1}^p \beta_i X_i,$$

donde  $\beta_i$  son los coeficientes del modelo.

- b) *Ajustamos o entrenamos* el modelo: utilizamos los datos para estimar los coeficientes.

## Métodos paramétricos

### Tensiones

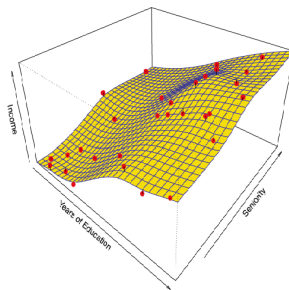
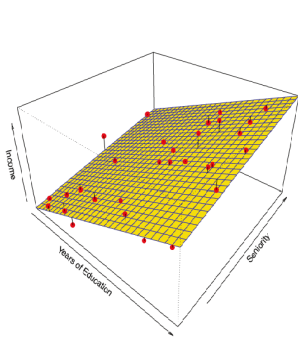
- Si  $f$  no es lineal, el ajuste será pobre.
- Si utilizamos modelos más flexibles (polinomios de orden superior por ejemplo) pagamos el precio de
  - estimar más parámetros;
  - posiblemente, sobreajustar los datos, *siguiendo al error demasiado cercanamente*.

¿Qué es la “ciencia de datos”?  
**Aprendizaje estadístico**  
La maldición de la dimensionalidad

Primeras nociones  
¿Por qué estimar  $f$ ?  
¿Cómo estimar  $f$ ?  
Aprendizaje supervisado versus no supervisado  
Regresión versus Clasificación  
Precisión del modelo  
El compromiso sesgo-varianza  
El contexto de clasificación

## Métodos paramétricos

Ajuste de los datos de ingreso versus escolaridad y antigüedad para  $n = 30$  personas (¿cuál es un mejor ajuste?).





## Métodos no-paramétricos

No se asume una forma particular para  $f$ .

- Las suposiciones son del tipo:  $f$  pertenece a la familia de las funciones continuas o a la de las funciones diferenciables, etc.
- Observamos que:
  - la ventaja consiste en la flexibilidad del método;
  - y su desventaja en que al no estimar un número finito de parámetros, más datos son necesarios para ajustar una función que pertenece a un espacio infinito-dimensional.

## Métodos no-paramétricos: estimador basado en núcleos

Más formalmente, dado el modelo

$$Y = f(\mathbf{x}) + \epsilon \quad (5)$$

y si asumimos que  $f$  es derivable de orden 2 con  $f''$  continua entonces un estimador consistente típico tiene la forma

$$\hat{f}(x) = \frac{\sum_{i=1}^n y_i K\left(\frac{x-x_i}{h_n}\right)}{\sum_{i=1}^n K\left(\frac{x-x_i}{h_n}\right)}$$

con  $K$  un núcleo y  $h_n$  es una sucesión de aberturas de ventana.

## ¿Qué método utilizamos?

Los métodos son **flexibles** o **restrictivos** según las familias de funciones  $f$  sean grandes o pequeñas.

$$\mathcal{F}_0 = \{f : f \text{ es lineal}\} \subset \cdots \subset \mathcal{F}_0 = \{f : f \text{ es continua}\}$$

Por ejemplo:

- una regresión lineal es un método poco flexible o muy **restrictivo**;
- una **regresión no paramétrica basada en núcleos** es muy flexible.

## ¿Qué metodo utilizamos?

### Compromiso entre precisión e interpretación

*¿Por qué razón optaríamos por un método de escasa flexibilidad?*

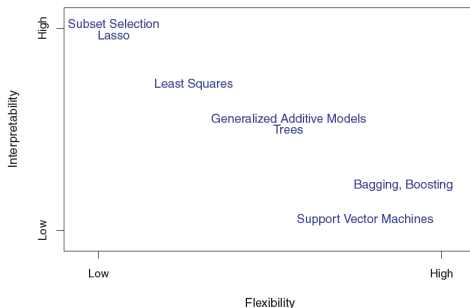
- Si estamos interesados en realizar una **inferencia**, un **modelo restrictivo** es siempre más simple de interpretar; por ejemplo un modelo lineal.
- Pero en cambio si sólo nos interesa **predecir** tal vez podamos optar por un **método más complejo**, menos interpretable pero con buena capacidad predictiva.

¿Qué es la “ciencia de datos”?  
**Aprendizaje estadístico**  
La maldición de la dimensionalidad

Primeras nociones  
¿Por qué estimar  $f$ ?  
¿Cómo estimar  $f$ ?  
Aprendizaje supervisado versus no supervisado  
Regresión versus Clasificación  
Precisión del modelo  
El compromiso sesgo-varianza  
El contexto de clasificación

## ¿Qué método utilizamos?

Diferentes métodos estadísticos se distribuyen en el “sistema de coordenadas” de interpretabilidad y flexibilidad.



**Figura 4:** “Distribución” de los métodos según el compromiso entre flexibilidad e interpretación. James et al. (2021)

## Aprendizaje supervisado versus no supervisado

La mayoría de los problemas estadísticos pueden agruparse en dos clasificaciones: *supervisado* y *no supervisado*.

En un problema **supervisado**;

- Contamos, para cada observación  $\mathbf{x}_i$  de un predictor con una observación de una **respuesta**,  $y_i$ ,  $i = 1, \dots, n$ .
- El objetivo es ajustar un modelo que relacione la respuesta a los predictores, en vistas de predecir y/o comprender esa relación (inferencia).
- **Regresión, regresión logística, modelos aditivos generalizados (GAM), boosting y support vector machines** operan en el dominio del análisis supervisado.

## Aprendizaje supervisado versus no supervisado

En un problema **no supervisado**

- contamos **sólo** con las observaciones del vector de **predictoras**  $\mathbf{x}_i$ ,  $i = 1, \dots, n$ , pero **no tenemos respuestas**,
- su denominación se debe a que carecemos de la respuesta (la variable  $Y$ ) que nos pueda guiar o supervisar nuestro análisis.
- **Cluster análisis** es una técnica típica dentro de las no supervisadas.

¿Qué es la “ciencia de datos”?  
**Aprendizaje estadístico**  
La maldición de la dimensionalidad

Primeras nociones  
¿Por qué estimar  $f$ ?  
¿Cómo estimar  $f$ ?  
Aprendizaje supervisado versus no supervisado  
**Regresión versus Clasificación**  
Precisión del modelo  
El compromiso sesgo-varianza  
El contexto de clasificación

## Regresión versus Clasificación

### Algunos tips

En “general”, si la respuesta es cuantitativa hablamos de **análisis de regresión** mientras que si es cualitativa hablamos de **clasificación**, aunque la distinción no puede ser establecida de forma tan tajante.



## Regresión versus Clasificación

### Algunos tips

- La **regresión por mínimos cuadrados** supone una variable respuesta cuantitativa.
- La **regresión logística** es utilizada para una variable respuesta cualitativa con dos respuestas posibles (binaria) y, al mismo tiempo, resuelve un problema de clasificación.
- **Vecinos más próximos** y **boosting** pueden ser utilizados tanto para variables respuesta cualitativas como cuantitativas.

## Precisión del modelo en el contexto de regresión

Queremos evaluar el **desempeño** de un método de aprendizaje estadístico **sobre un conjunto de datos dados**, por ende necesitamos una medida de cuán bien sus **predicciones ajustan** a los **datos observados**.

Una posible medida es el **error cuadrático medio (estimado)**

$$\widehat{\text{ECM}} = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{f}(x_i))^2 \quad (6)$$

donde  $\hat{f}$  es el estimador de  $f$ .

## Precisión del modelo en el contexto de regresión

Es importante distinguir dos tipos de problemas:

- a) Si asumimos distribución normal, independencia y que el modelo es lineal  $f$ , **el problema de optimalidad de la selección de  $f$  está resuelto** y hay unicidad de esa solución.
- b) El escenario al que se refiere James et al. (2021) es **cuando no hay cumplimiento claro de supuestos** y entonces aparecen varios métodos competidores de aprendizaje estadístico: regresión lineal simple, splines, regresión no paramétrica para el mismo conjunto de datos.

En ambos contextos seleccionamos aquel  $\hat{f}$  que minimice  $\widehat{ECM}^{11}$  pero las situaciones son bien diferentes.

---

<sup>11</sup>Por simplicidad escribiremos ECM

## Precisión del modelo en el contexto de regresión

Observemos que:

- El ECM está calculado sobre los datos dados, a los que llamaremos **datos de entrenamiento**.
- Y estamos interesados en la precisión de las predicciones que obtenemos a partir de los **datos test** todavía no vistos.

## Datos de entrenamiento y datos test

Más generalmente, supongamos que contamos con una colección

$$C = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)\}$$

de  $n$  datos de entrenamiento.

En base a  $C$  se construyó  $\hat{f}_C$  y el  $\text{ECM}(\hat{f}_C)$  mide la calidad del ajuste.

Un buen método de aprendizaje estadístico es aquel que, además de minimizar el  $\text{ECM}(\hat{f}_C)$ , frente a un nuevo dato test  $(x_0, y_0) \notin C$  satisfaga

$$y_0 \approx \hat{f}_C(x_0).$$

## Datos de entrenamiento y datos test

Y si contamos con una colección grande

$$C_N = \{(\mathbf{x}_1^0, y_1^0), \dots, (\mathbf{x}_m^0, y_m^0)\}$$

de “nuevos” datos interesa el procedimiento  $\hat{f}_C$  que minimice tanto (6) como

$$\frac{1}{m} \sum_{i=1}^m (y_i^0 - \hat{f}_C(x_i^0))^2. \quad (7)$$

A esta expresión le denominamos ECM-test o error test.

## Datos de entrenamiento y datos test: ejemplo

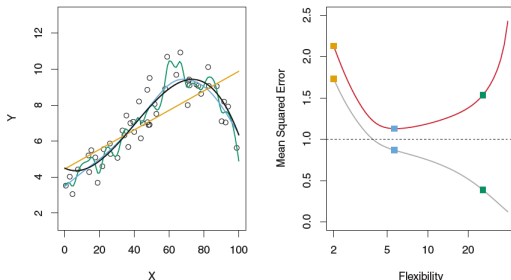


Figura 5: Desempeño de tres métodos de estimación de  $f$ .

Se generan observaciones a partir  $f$  (conocida) y se comparan tres métodos de estimación de  $f$ . En el panel izquierdo se representan: estimación por regresión lineal (línea naranja) y por splines (líneas azul y verde). La línea negra sólida corresponde a la curva de  $f$ . En el panel derecho se comparan los valores de ECM para los datos de entrenamiento (curva gris), para los datos test (curva roja) y el mínimo ECM-test sobre todos los métodos (línea quebrada). James et al. (2021)

## Datos de entrenamiento y datos test: ejemplo

**Observaciones.** En relación al panel izquierdo de la Figura 5.

- Las líneas **naranja**, **azul** y **verde** representan tres estimadores para  $f$ , utilizando métodos con niveles de flexibilidad creciente.
- La **línea naranja** indica un **ajuste por regresión lineal**, relativamente **inflexible**. La **azul** y **verde** fueron producidos por splines con diferentes niveles de suavizado.
- A medida que aumenta el suavizado las curvas de los estimadores se aproximan mejor a los datos.
- La **curva verde** es la **más flexible** y se parea con los datos muy bien pero ajusta pobremente a  $f$  por su alta variabilidad.



## Datos de entrenamiento y datos test: ejemplo

**Observaciones.** Panel derecho de la Figura 5.

- La línea gris muestra el ECM basado en los datos de entrenamiento como una función de la flexibilidad. Los cuadrados **naranja**, **azul** y **verde** indican los valores de ECM asociados con las correspondientes curvas del panel izquierdo.
- Utilizaremos los términos “**grados de libertad**” para indicar, provisoriamente, “**flexibilidad de una curva**”. Cuanto mayor ondulaciones tenga una curva más grados de libertad posee.

## Datos de entrenamiento y datos test: ejemplo

**Observaciones.** Panel derecho de la Figura 5.

- El ECM-entrenamiento decrece monótonamente cuando la flexibilidad aumenta.
- Como  $f$  no es lineal, el “ajuste naranja” es pobre mientras que la curva verde muestra el más bajo de los ECM-entrenamiento por su alto grado de libertad.

## Datos de entrenamiento y datos test: ejemplo

**Observaciones.** En relación al panel derecho de la Figura 5.

- La línea roja muestra el ECM basado en los datos test.
- Similar a ECM-entrenamiento, ECM-test decrece inicialmente cuando aumenta la flexibilidad pero a partir de cierto punto comienza a aumentar.
- Así, las curvas naranja y verdes tienen peor desempeño que la azul (mayor ECM-test).
- La línea horizontal quebrada indica  $\text{VAR}(\epsilon)$ , el componente irreducible de (2) y por ende el mínimo ECM-test alcanzable entre todos los métodos posibles.

## Más allá del ejemplo

### Dimensiones generales del entrenamiento estadístico

- El decrecimiento ECM-entrenamiento a medida que la flexibilidad aumenta y la forma de U de ECM-test versus flexibilidad es un fenómeno general de los **métodos de aprendizaje estadístico** (mae).
- Cuando la flexibilidad aumenta ECM-entrenamiento decrecerá pero puede que ECM-test no.

## Más allá del ejemplo

### Dimensiones generales del entrenamiento estadístico

Cuando un método produce un ECM-entrenamiento muy pequeño y un ECM-test grande estamos frente a un **sobreajuste de los datos** que:

- se debe a que nuestro mae trabaja demasiado exhaustivamente para encontrar patrones en los datos de entrenamiento y puede que halle algunos patrones que **se deben sólo al fenómeno aleatorio** y no a las propiedades verdaderas de  $f$ ;
- provoca un ECM-test grande, como efecto de que **los patrones** que el mae halla en los datos de entrenamiento simplemente **no existen** en los datos test.

## Más allá del ejemplo

### Dimensiones generales del entrenamiento estadístico

- En general

ECM-entrenamiento  $<$  ECM-test.

.

- Cuando hablamos de **sobreajuste** de un mae es porque hay otro mae menos flexible que ha producido ya un ECM-test menor.

## Más allá del ejemplo

- En muchos contextos podemos contar con observaciones test, pero en general no
- ¿Y si sólo contamos con los datos de entrenamiento? Hay que construir una estrategia de expresión máxima de los datos de tal modo que el método aprenda (**convalidación o validación cruzada, por ejemplo**).

## Descomposición sesgo-varianza

Se puede demostrar que, para el modelo (1), en condiciones muy generales

$$\begin{aligned}\text{ECM} \left( \hat{f}(x_0) \right) &=: E \left( y_0 - \hat{f}(x_0) \right)^2 \\ &= \text{VAR} \left( \hat{f}(x_0) \right) + \left[ \text{SES} \left( \hat{f}(x_0) \right) \right]^2 + \text{VAR}(\epsilon) \quad (8)\end{aligned}$$

donde

$$\text{SES} \left( \hat{f}(x_0) \right) = E(\hat{f}(x_0)) - f(x_0)$$

es el sesgo del estimador.



## Descomposición sesgo-varianza

Observaciones importantes:

- El ECM global se obtiene promediando (computo una integral o una serie)  $\text{ECM} \left( \hat{f}(x_0) \right)$  sobre todos los puntos  $x_0$ .

- $\forall x_0$ :

$$\text{ECM} \left( \hat{f}(x_0) \right) \geq \text{VAR}(\epsilon)$$

- La igualdad (8) nos indica que un buen mae debe tener bajos sesgo y varianza.
- Un buen método debe minimizar el El ECM global.

## Compromiso sesgo-varianza

### Componente varianza

- Elegido el mae, diferentes estras de datos de entrenamiento producen distintos  $\hat{f}$ . Un buen mae no debería variar demasiado entre las diferentes muestras. Esta es la noción de **varianza** del método.
- En general los mae muy flexibles tienen alta **varianza**.

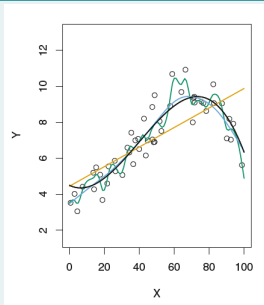
Recordemos el panel izquierdo de la Figura 5, ahora Figura 6.

¿Qué es la “ciencia de datos”?  
Aprendizaje estadístico  
La maldición de la dimensionalidad

Primeras nociones  
¿Por qué estimar  $f$ ?  
¿Cómo estimar  $f$ ?  
Aprendizaje supervisado versus no supervisado  
Regresión versus Clasificación  
Precisión del modelo  
El compromiso sesgo-varianza  
El contexto de clasificación

## Compromiso sesgo-varianza

### Componente varianza



**Figura 6:** Desempeño de tres métodos de estimación de  $f$ . La curva **verde** (**naranja**) representa un estimador **muy variable** ya que sigue a los datos (**poco flexible y con baja varianza**).

## Compromiso sesgo-varianza

### Compromiso sesgo–varianza

Figura 7:

- Cuando  $f$  no es lineal, en un ajuste lineal (baja flexibilidad) el aporte al **MSE** lo realiza el **sesgo** y el error es grande. En cambio si la flexibilidad es alta el **MSE** es bajo y el **sesgo** es casi nulo.

Figura 8:

- Cuando  $f$  es casi lineal un ajuste lineal (baja flexibilidad) minimiza **MSE**. Un ajuste no lineal (muy flexible) provoca un **sesgo** casi nulo pero un alto **MSE** el cual se explica por la **varianza**.

¿Qué es la “ciencia de datos”?  
Aprendizaje estadístico  
La maldición de la dimensionalidad

Primeras nociones

¿Por qué estimar  $f$ ?

¿Cómo estimar  $f$ ?

Aprendizaje supervisado versus no supervisado

Regresión versus Clasificación

Precisión del modelo

El compromiso sesgo-varianza

El contexto de clasificación

## Compromiso sesgo-varianza

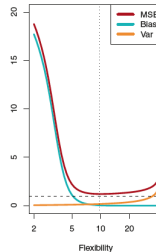
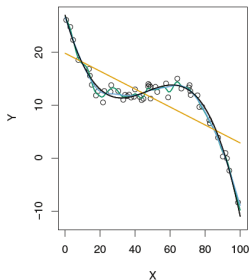


Figura 7: Contribución del sesgo y la varianza cuando  $f$  no es lineal

¿Qué es la “ciencia de datos”?  
**Aprendizaje estadístico**  
La maldición de la dimensionalidad

Primeras nociones

¿Por qué estimar  $f$ ?

¿Cómo estimar  $f$ ?

Aprendizaje supervisado versus no supervisado

Regresión versus Clasificación

Precisión del modelo

**El compromiso sesgo-varianza**

El contexto de clasificación

## Compromiso sesgo-varianza

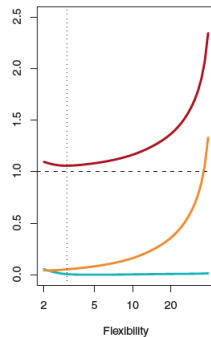
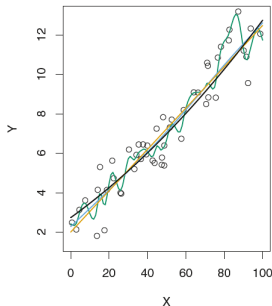


Figura 8: Contribución del sesgo y la varianza cuando  $f$  es casi lineal

## Clasificación

Como antes consideremos  $n$  observaciones de  $(\mathbf{x}, Y)$ ,  $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)$ , donde la variable respuesta  $Y$  es **cualitativa**.

Ejemplos típicos:

- Una persona arriba a un centro de emergencia con un conjunto de síntomas (predictoras) que podrían ser atribuidos a tres condiciones médicas (respuesta).
- Una línea de servicio de un banco determina cuando una transacción es fraudulenta de acuerdo a dirección IP del usuario, historia pasada de las transacciones, etc.
- De acuerdo al largo y ancho del sépalos y del pétalo se clasifica la especie de una flor Iris.

## Clasificación

Nuestro estimador  $\hat{f}$  es un **clasificador** (regla).

### Tasa de error de entrenamiento

- La precisión de  $\hat{f}$  es la tasa de error de entrenamiento: la proporción de errores que se producen cuando aplicamos  $\hat{f}$  a los datos de entrenamiento; i.e.

$$\frac{1}{n} \sum_{i=1}^n I(y_i \neq \hat{y}_i)$$

donde  $\hat{y}_i$  es la etiqueta de clase predicha para la  $i$ -ésima observación utilizando  $\hat{f}$ .



## Clasificación

Nuestro estimador  $\hat{f}$  es un **clasificador** (regla).

### Tasa de error test

- La *tasa de error test* se define de manera análoga para un conjunto de observaciones test. Más específicamente, dado  $(\mathbf{x}_0, y_0)$ , esa tasa es

$$\text{promedio}(I(y_0 \neq \hat{y}_0)).$$

- *Un buen clasificador es aquel que hace mínima esta tasa de error test.*

## Clasificador de Bayes

La *tasa de error test* es minimizada, en promedio, por el **clasificador de Bayes** que *asigna a cada observación la clase más probable dado los valores del predictor*.

Formalmente, dado  $\mathbf{x} = \mathbf{x}_0$  le asignamos  $Y = j_0$  si

$$P(Y = j_0 | \mathbf{X} = \mathbf{x}_0) = \max_j P(Y = j | \mathbf{X} = \mathbf{x}_0).$$

Así, este clasificador es el mejor u óptimo.

## Clasificador de Bayes

### Clasificador de Bayes con dos clases 1 y 2

Dado  $\mathbf{x} = \mathbf{x}_0$ , el clasificador de Bayes corresponde a predecir la clase 1 si  $P(Y = 1|\mathbf{x} = \mathbf{x}_0) > \frac{1}{2}$  y, la clase 2, si es de otro modo.

### Frontera de decisión de Bayes

se define como el conjunto de puntos

$$\left\{ \mathbf{x} : P(Y = 1|\mathbf{x} = \mathbf{x}) = P(Y = 0|\mathbf{x} = \mathbf{x}) = \frac{1}{2} \right\}.$$

## $k$ -vecinos más próximos

- El clasificador de Bayes, que es el optimal, supone conocidas las probabilidades de clasificación

$$P(Y = j_0 | \mathbf{X} = \mathbf{x}_0), \forall j, \forall \mathbf{x}_0.$$

- Pero para un conjunto de datos reales (que no sean simulados de una distribución conocida) estas probabilidades no se conocen y se deben estimar.

## $k$ -vecinos más próximos

El estimador de “ $K$ —vecinos más próximos” (KNN) procede por estimación de dichas probabilidades.

Construyamos el estimador por pasos:

**Paso 1** Dado un entero positivo  $K$  y una observación test  $\mathbf{x}_0$ , el KNN identifica los  $K$  puntos en la muestra de entrenamiento más cercanos (utilizando alguna métrica) a  $\mathbf{x}_0$ . Denotamos con  $\mathcal{N}_{\mathbf{x}_0}$  al conjunto de los vecinos.

## $k$ -vecinos más próximos

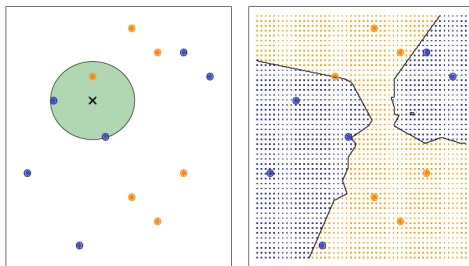
**Paso 2** Estimemos la probabilidad condicional

$\pi_{j|\mathbf{x}_0} = P(Y = j_0 | \mathbf{X} = \mathbf{x}_0)$  como

$$\hat{\pi}_{j|\mathbf{x}_0} = \frac{1}{K} \sum_{i \in \mathcal{N}_{\mathbf{x}_0}} I(y_i = j)$$

**Paso 3** El clasificador KNN aplica la regla de Bayes: a la observación  $\mathbf{x}_0$  le asigna la clase  $j_0$  que maximiza  $\hat{\pi}_{j|\mathbf{x}_0}$  en la colección de todas las clases.

## $k$ -vecinos más próximos



**Figura 9:** KNN con  $K = 3$  vecinos más próximos<sup>12</sup>. En el panel izquierdo se exhibe la nueva observación  $x$  clasificada como azul. En el panel derecho se representa la frontera de decisión y las clases.

<sup>12</sup>Tomada de James et al. (2021)

## Introducción: datos de alta dimensión

Son consecuencia del registro simultáneo de miles (o millones) de características sobre cada objeto o individuo

- Las biotecnologías asociadas a la genética como los microarreglos.
- Imágenes y videos: imágenes médicas, astrofísicas, de vigilancia, etc.
- Datos de preferencias de consumo: programas de fidelización, redes sociales, etc.
- Datos provistos por los negocios: compañías de logística y transporte, de seguros, de la industria financiera, etc.
- Datos colaborativos: los datos surgen de los registros de miles de usuarios a través de celulares. Ejemplo: “eBirth”.



## ¿Bendición o maldición?

¿No es acaso, para la estadística, una bendición poder registrar miles de variables para un individuo?

La respuesta es... ¡depende!

En estadística de alta dimensión separar en general la señal del ruido es casi imposible

## La maldición de la dimensionalidad

El impacto de una dimensión alta puede ser enorme y paradójal.  
En  $\mathbb{R}^p$  con  $p$  grande

- Los puntos son aislados
- Las acumulaciones de pequeñas fluctuaciones en muchas direcciones diferentes pueden producir una gran fluctuación global
- Un evento que es acumulación de eventos raros puede no ser raro
- Los cálculos numéricos y las optimizaciones pueden tener alto costo

## Los espacios son vastos

### Modelo de Regresión

$$Y_i = f(\mathbf{X}^{(i)}) + \epsilon_i \quad (9)$$

con

- $\mathbf{X}^{(i)}$ ,  $i = 1, \dots, n$  observaciones i.i.d de un vector  
 $\mathbf{X} = (X_1, \dots, X_p)'$  de dimensión  $p$
- $\epsilon_1, \dots, \epsilon_n$  son i.i.d y están centradas.

Estimamos  $f(\mathbf{x})$  utilizando *información local*, por ejemplo con  $k$ -vecinos más próximos.

Si  $p$  es grande ¿cuántas observaciones hay cercanas a  $\mathbf{X} = \mathbf{x}$ ?

## Los espacios son vastos

Una pregunta previa a la anterior:

¿Cómo se comportan las distancias  $\|\mathbf{x}^{(i)} - \mathbf{x}^{(j)}\|^2$ ?

Para simplificar  $X_1, \dots, X_p$  son i.i.d con  $X_1 \sim U([0, 1])$

↓

$$\mathbf{X} = (X_1, \dots, X_p)' \sim U([0, 1]^p)$$

Las Figuras 10 y 11 muestran los cuatro histogramas de

$$\|\mathbf{x}^{(i)} - \mathbf{x}^{(j)}\|, 1 \leq i < j \leq n = 100$$

en cuatro escenarios (dimensiones):

$$p = 2, 10, 100, 1000.$$

¿Qué es la “ciencia de datos”?  
Aprendizaje estadístico  
La maldición de la dimensionalidad

Introducción

Perdidos en la inmensidad de espacios de alta dimensión  
Regresión lineal en espacios de alta dimensión  
Complejidad computacional  
Cambio de paradigma

## Los espacios son vastos

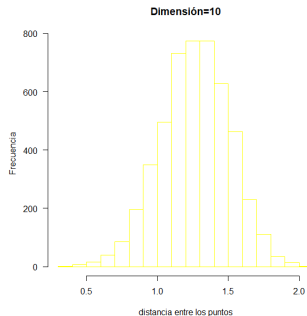
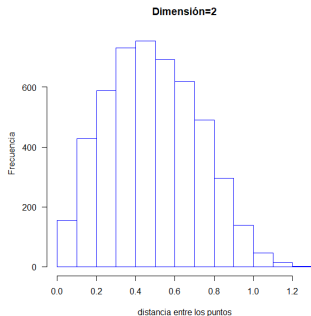


Figura 10: Distribuciones empíricas de las distancias,  $p = 2, 10$ .

¿Qué es la “ciencia de datos”?  
Aprendizaje estadístico  
La maldición de la dimensionalidad

Introducción

Perdidos en la inmensidad de espacios de alta dimensión  
Regresión lineal en espacios de alta dimensión  
Complejidad computacional  
Cambio de paradigma

## Los espacios son vastos

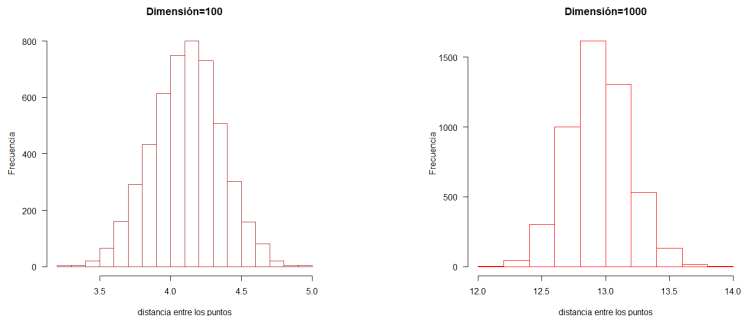


Figura 11: Distribuciones empíricas de las distancias,  $p = 100, 1000$ .

## Los espacios son vastos

De los resultados empíricos

- ❑ La distancia entre los puntos se incrementa con  $p$ .
- ❑ Si  $p$  es grande **todos los puntos están a una distancia similar**.
- ❑ Notar que la DCM y el DE <sup>13</sup> satisfacen

$$\square \text{ DMC} = E \left\| \mathbf{x}^{(i)} - \mathbf{x}^{(j)} \right\|^2 = p/6$$

$$\square \text{ SD} = \left( \sum_{k=1}^p \text{VAR} \left[ \left( \mathbf{x}^{(i)} - \mathbf{x}^{(j)} \right)^2 \right] \right)^{1/2} \approx 0.2\sqrt{p}$$

$$\square \text{ SD/DMC} \sim p^{-1/2}.$$

<sup>13</sup> Distancia Media Cuadrática y Desvío Estándar, respectivamente

## Los espacios son vastos

Dado  $\mathbf{x} \in [0, 1]^p$ ,

$$\text{¿}n : \exists \mathbf{X}^{(i)} : \|\mathbf{X}^{(i)} - \mathbf{x}\| \leq 1 \text{ ?}$$

El volumen de la bola  $p$ -dimensional de radio  $r$  es

$$V_p(r) = \frac{\pi^{p/2}}{\Gamma(\frac{p}{2} + 1)} r^p \sim \left( \frac{2\pi e r^2}{p} \right)^{p/2} (p\pi)^{-1/2} \quad (10)$$

$p \rightarrow \infty$

Si

$$[0, 1]^p \subseteq \bigcup_{i=1}^n B(\mathbf{X}^{(i)}, 1)$$

entonces



## Los espacios son vastos

Volviendo a la pregunta, dado  $\mathbf{x} \in [0, 1]^p$ ,

$$\text{¿} n : \exists \mathbf{X}^{(i)} : \left\| \mathbf{X}^{(i)} - \mathbf{x} \right\| \leq 1 ?$$

- Si  $p = 20$  entonces  $n = 39$
- Si  $p = 50$  entonces  $n = 5.7 \times 10^{12}$
- Si  $p = 200$  entonces

$n$  es mayor que el número total de partículas observables en el universo

## Los espacios son vastos

De (10) se deduce que:

- $V_p(r) \downarrow 0$  con una tasa más rápida que una exponencial si  $p \uparrow \infty$
- por ejemplo, para  $p = 20$  el volumen de la bola unidad es (casi) cero!
- si  $C_p(r) = \{\mathbf{x} \in [0, 1]^p : 0.99r \leq \mathbf{x} < r\}$  entonces

$$\frac{\text{volumen}(C_p(r))}{\text{volumen}(B_p(r))} \sim 1$$

exponencialmente cuando  $p \rightarrow \infty$ .

- si  $p$  es grande, “la mayoría” de los puntos de la bola  $B_p(r)$  están más cerca de su frontera que de su centro.

## Regresión Lineal en alta dimensión

Consideremos el siguiente modelo de regresión

$$Y_i = \langle \mathbf{X}^{(i)}, \beta^* \rangle + \epsilon_i, i = 1, \dots, n,$$

con  $\mathbf{X}$  la matriz de datos, de la que asumimos que sus columnas son ortonormales. Entonces

$$\mathbb{E} \left[ \left\| \hat{\beta} - \beta^* \right\|^2 \right] = p\sigma^2$$

con lo cual **si la dimensión  $p$  es grande el error es grande.**

## Regresión Lineal en alta dimensión

**Ejemplo:** Consideremos el modelo

$$Y_i = f_{\beta^*}(i/n) + \epsilon_i, \quad i = 1, \dots, n$$

siendo

$$f_{\beta^*}(t) = \sum_{j=1}^p \beta_j^* \cos(\pi j t)$$

y donde  $\epsilon_1, \dots, \epsilon_n$  es una muestra aleatoria de una normal estándar.

En la Figura 12 se representan 100 observaciones del modelo y las curvas correspondientes a la función de regresión y al estimador de mínimos cuadrados en cuatro escenarios:  $p = 10, 20, 50, 100$ . Los  $\beta_j^*$  son constantes muestreadas de una  $N(0, j^{-4})$ .

¿Qué es la “ciencia de datos”?  
Aprendizaje estadístico  
La maldición de la dimensionalidad

Introducción  
Perdidos en la inmensidad de espacios de alta dimensión  
Regresión lineal en espacios de alta dimensión  
Complejidad computacional  
Cambio de paradigma

## Regresión Lineal en alta dimensión

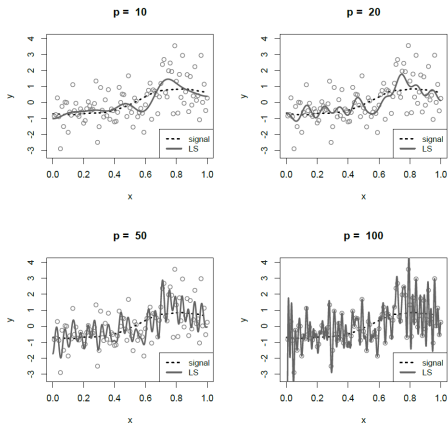


Figura 12: Estimador de mínimos cuadrados  $f_{\hat{\beta}}$ , con  $n = 100$  y  $p = 10, 20, 50, 100$ . Los puntos son

## Complejidad computacional

En un contexto de alta dimensión surgen

- una gigantesca intensidad del cómputo numérico
- e incluso la posibilidad de exceder las fuentes o recursos (memoria) computacionales

## Cambio de paradigma

En estadística clásica hay una rica teoría para analizar los datos cuando  $n$  es grande y  $p$  pequeño.

Pero en muchos campos científicos esto no ocurre y, en cambio,

- hay un número grande de parámetros  $p$ .
- puede haber tamaños de muestras del mismo tamaño que  $p$  o incluso de tamaño menor.

En consecuencia, en este nuevo contexto,

la teoría asintótica “ $p$  fijo y  $n \rightarrow \infty$ ” no tiene sentido.