

# Big Data, Machine Learning and Econometrics

Un tour por el mundo del Machine Learning

Gabriel Martos

gmartos@utdt.edu



UNIVERSIDAD  
TORCUATO DI TELLA

# Agenda

- 1 Organización del curso
- 2 Contextualización del curso
- 3  $\hat{\beta}$  vs  $\hat{y}$
- 4 Breve intro a R en R

**Objetivo:** Introducir los modelos de aprendizaje supervisado y aprender a entrenar los mismos en R (casos estudio: Economía, Finanzas y Negocios).

## A– Introducción:

- ▶ Contextualización del curso e introducción a R.
- ▶ Breve intro a métodos numéricos de optimización en ML.
- ▶ Minimización del riesgo empírico y descomposición del riesgo predictivo.

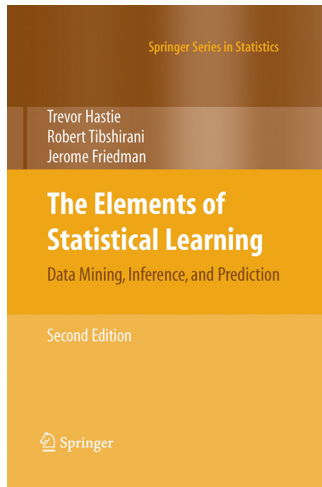
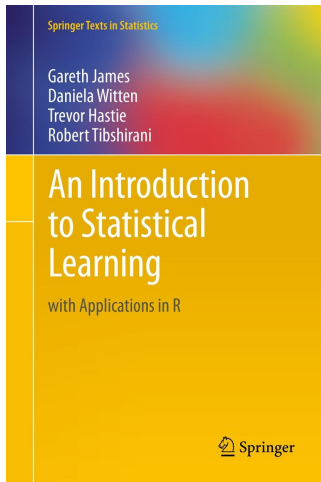
## B– Clasificación y Regresión:

- ▶ Validación cruzada y modelo de K-vecinos.
- ▶ Modelos Aditivos, Selección y Regularización (Ridge, Lasso, ENets).
- ▶ Árboles y ensamble de modelos (Bagging, Random Forest, Boosting).
- ▶ Métodos de Kernel y Máquinas de Vector Soporte (SVM).
- ▶ Introducción a la modelización con Redes Neuronales.

## C– Modelos predictivos para series temporales:

- ▶ Singular Spectrum Analysis y Procesos Gaussianos para modelar y predecir series temporales (VC para modelos de series).

# Bibliografía



(Indicada al principio de cada tema)

# Evaluación

- Trabajo práctico individual.
- Fecha límite de entrega: Febrero del 2025.
- Aspectos a evaluar:
  - ▶ Motivación y contextualización del problema predictivo.
  - ▶ Revisión bibliográfica relevante sobre el problema.
  - ▶ Originalidad de la presentación.
  - ▶ Contrastación de métodos y comparativa con modelos econométricos.
  - ▶ Claridad y profundidad en la presentación de la metodología que utilizaste para resolver el problema. Claridad en la exposición de los principales resultados obtenidos (ver el TP de referencia).

Algunos repositorios que puedes comenzar a explorar en búsqueda de datos / problemas de aprendizaje supervisado:

- [AcademicTorrents](#).
- [Open Data \(PLOS\)](#).
- [Datasetsearch](#).
- [Kaggle](#).
- [UCI Machine Learning repository](#).
- [OpenML](#).
- [DrivenData](#).
- [CrowdAnalytix](#).
- ...

# Agenda

- 1 Organización del curso
- 2 Contextualización del curso
- 3  $\hat{\beta}$  vs  $\hat{y}$
- 4 Breve intro a R en R

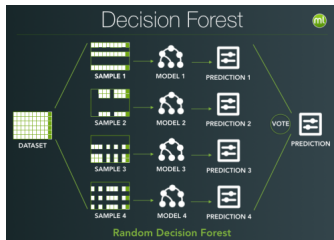
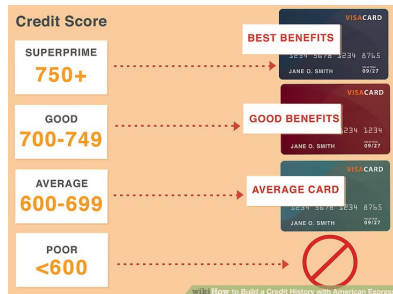
**Fact:** Las TIC's producen una cantidad ingente de DATOS.

- *Bigdata*: Conjunto de datos grande y complejo.
- ¿Cuanto “big” es el Bigdata?:
  - ▶ Empresas Tecnológicas (Google) almacena las consultas en su buscador de manera diaria (BD del orden de peta-bytes– $10^6$ GB– diarios!)
  - ▶ Bancos: Registran “eventos” diarios de clientes por distintos canales (cajero, homebanking, canales móvil, call-center, etc)... solo estas transacciones representan del orden de varios TB ( $10^3$ GB) por día.
  - ▶ **Organismos públicos** registran pagos de impuestos, trámites, consultas y reclamaciones, *vacunaciones y desplazamiento*, seguridad, etc.
- Necesitamos adaptar las **herramientas de análisis** de datos tradicionales para transformar los Bigdatos en información útil.
  - ▶ La maldición de la dimensión ( $p \gg 0$ ).
  - ▶ *Fatdata* ó  $p > n$ : Más regresores que observaciones.
  - ▶ NO linealidad en las relaciones, el problema de la multicausalidad, etc.
- ML como complemento del toolkit “estadísticonométrico”.



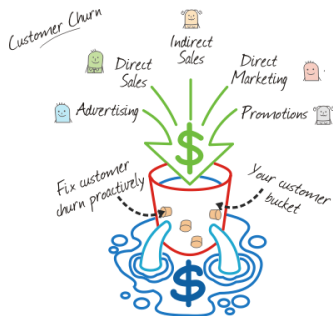
# Aplicaciones en Finanzas

- Instituciones financieras recopilan información de sus clientes.
- Modelos de “Score Crediticio”.
- La rentabilidad del negocio depende de manera crucial de la medición acertada del riesgo.



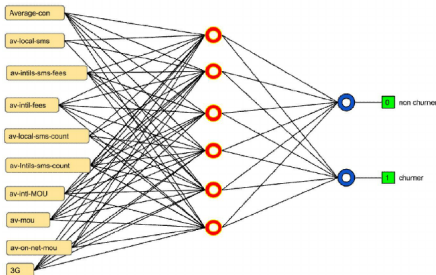
- Modelizaremos el riesgo crediticio con: Bagging, Random Forest, y Boosting.
- Capturamos efectos no lineales.
- Reducimos variabilidad ensamblando.

# Aplicaciones en Marketing (churn)



- Modelizamos la probabilidad de fuga.
- Para las firmas el output de estos modelos es esencial para planificar su oferta de bienes y servicios; como así también para diseñar estrategias de retención de clientes (promociones).

Redes neuronales para modelizar la fuga de clientes de telefonía.



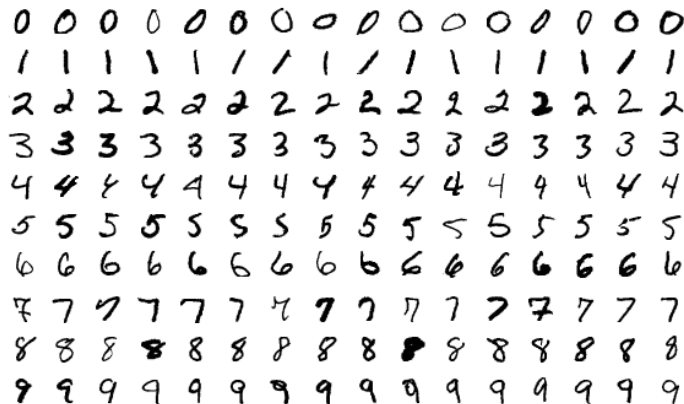
# Aplicaciones en Transporte y Energía

Objetivo: Modelizar la demanda para gestionar una producción y tarificación adecuada a lo largo del tiempo (cada hora, día, semana, mes, año, etc).



- Urban mobility lab (MIT): <https://mobility.mit.edu/machine-learning>

# Aplicaciones en Ciencia y Tecnología



Reconocimiento de imágenes para el diagnóstico médico, la seguridad informática, la automatización de procesos productivos, etc.

# Otras aplicaciones en negocios, economía

- Customer Segmentation: Detectar nichos de clientes con características específicas para orientar acciones comerciales.
- Pricing: Modelizar la demanda de productos para gestionar una producción y tarificación adecuada a lo largo del tiempo.
- Recommendation Engines: Predecir de un conjunto muy grande de datos las preferencias/gusto de los clientes:

`https://www.netflixprize.com/`

- ...

# Agenda

- 1 Organización del curso
- 2 Contextualización del curso
- 3  $\hat{\beta}$  vs  $\hat{y}$
- 4 Breve intro a R en R

*... my advice to grads (in Economics) is “go to the computer science department and take a course in machine learning.”*

H. Varian (Chief economist at Google)

- Varian, H. R. (2014): “*Big data: New tricks for econometrics*”. Journal of Economic Perspectives, 28(2), 3-28.

Estadística y Econometría están orientadas a:

- Hacer *Inferencia* apoyándose en **modelos probabilísticos**.
- Los modelos sirven para describir y cuantificar como se relación las variables ( $\hat{\beta}$ ). Nos interesan los parámetros del modelo per-se.
- En menor medida a hacer predicciones con el modelo.



## Estadística y Econometría están orientadas a:

- Hacer *Inferencia* apoyándose en **modelos probabilísticos**.
- Los modelos sirven para describir y cuantificar como se relación las variables ( $\hat{\beta}$ ). Nos interesan los parámetros del modelo per-se.
- En menor medida a hacer predicciones con el modelo.

## Machine Learning orientado a:

- **Predicción** de las variables de interés ( $\hat{y}$ ). Los parámetros de los modelos interesan poco (son sólo un medio para obtener  $\hat{y}$ ).
- Filosofía del “*dejar que los datos hablen por si mismos*” (modelos no-paramétricos y con escasa/nula estructura probabilística).

# ¿Qué pueden aprender los econometristas del ML?

- Overfitting (explicitar el trade-off sesgo y varianza de cada modelo).
- Técnicas de Validación Cruzada:
  - ▶ Aprender valores adecuados de ciertos (hiper/meta)parámetros.
  - ▶ Estimar la performance out-of-sample.
- Regularización: Ridge, Lasso “and friends” (HD data).
- Ensamble de modelos (otro mecanismo de regularización).
- Modelos no lineales: Árboles, SVM, Redes Neuronales, etc.
- Herramientas computacionales avanzadas y escalabilidad de los métodos de análisis al contexto del Bigdata.
- Herramientas para explorar nuevos dominios: Análisis semántico, de imágenes, de redes sociales, análisis de sentimientos, etc.

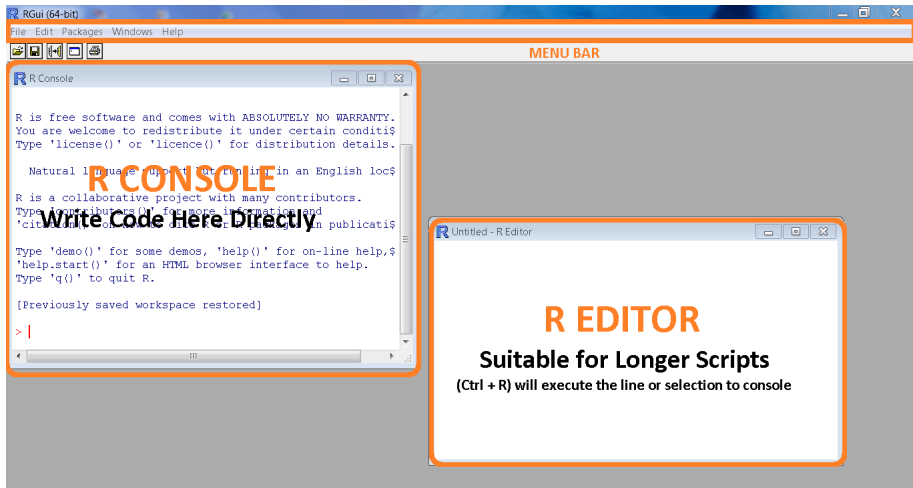
# Algunos economistas que trabajan con modelos de ML:

- Susan Athey (Stanford Graduate School of Business).
- Guido Imbens (Stanford Graduate School of Business).
- Jann Spiess (Harvard).
- Andrew Gelman (Universidad de Columbia).
- Interesante para seguir leyendo:
  - ▶ *Why a leading economist is embracing machine learning.*
  - ▶ *Machine Learning Methods Economists Should Know About*

# Agenda

- 1 Organización del curso
- 2 Contextualización del curso
- 3  $\hat{\beta}$  vs  $\hat{y}$
- 4 Breve intro a R en R

- En el curso vamos a utilizar un enfoque “metodológico–práctico” para presentar cada uno de los modelos y temas.
- Para correr estos modelos, vamos a utilizar R.
- R es un lenguaje GRATUITO de programación ampliamente utilizado (al igual que Python) en Ciencia de Datos. Links para descargarlo:
  - ▶ [Baja primero R de su web oficial.](#)
  - ▶ Aconsejable descargar también [RStudio](#) (yo lo voy a utilizar en clase).
- Vamos a invertir el resto de la clase en aprender las instrucciones y comandos básicos para:
  - ▶ Crear, cargar y manipular datos (esto último es lo más trabajoso).
  - ▶ Utilizar estructuras de control que luego serán de muchísima utilidad (por ejemplo para implementar validación cruzada con un modelo).
- Una vez instales ambos programas verás lo siguiente...



The screenshot shows the RStudio environment with the following components:

- Code Area:** Contains R code for a Shiny application. The code includes `library(shiny)`, `shinyUI(pagewithsidebar(...))`, and `mainPanel(...)` with a `selectInput` for chart type.
- Workspace Area:** Displays the current workspace contents. It shows a data frame `bse` with 3924 observations and 7 variables. The variables are `x` (Date[3924]) and `y` (numeric[3924]).
- Console Area:** Shows the execution of R commands. The commands include `view(bse)`, `view(bse)`, `bseDate <- as.Date(bse$date, format="%Y-%m-%d")`, and several `plot` calls for `bse` data, including `plot(x=bseDate, y=bse$open, type="l", main="BSE Data", col="blue", xlab="Periods", ylab="Index", lwd=3)`. An error message "Error in plot.new() : figure margins too large" is visible.
- Plot, Help and Package Area:** Displays a line plot titled "BSE Data". The x-axis is labeled "Periods" and the y-axis is labeled "Index". The plot shows a blue line representing the BSE data over time, with a peak around 2000.