

# ¿Qué es Machine Learning?

## Fundamentos conceptuales

German Rosati  
german.rosati@gmail.com

UNSAM - CONICET

22 de octubre de 2020

# Preguntas

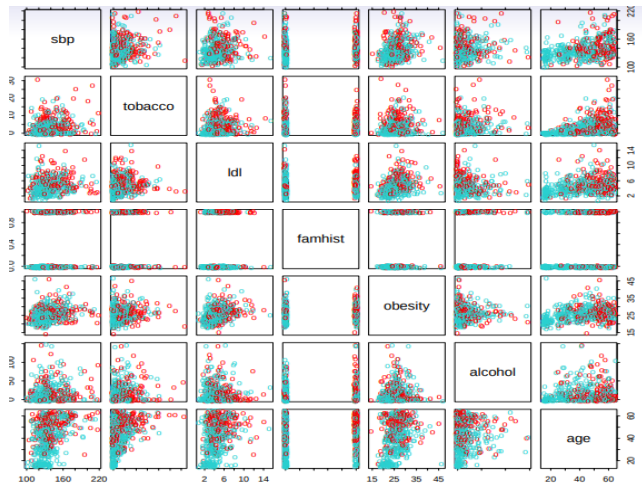
- ¿Podría una computadora ir más allá de “lo que sea que sepamos decirle que haga” y realmente “aprender” por su cuenta como realizar una determinada tarea?
- ¿Podría ser posible el aprendizaje automático de estas reglas a partir de los datos?

## ¿Qué problemas se pueden abordar?

- Riesgo de cáncer prostático
- Predicción de ataques cardíacos en función de variables demográficas, orgánicas, médicas, etc.
- Generar un detector de spam automático para tu cuenta de correo
- Reconocer dígitos en números manuscritos
- Clasificar muestras de tejido en diferentes tipos de cáncer basado en ciertos perfiles genéticos
- Establecer relaciones entre el ingreso y variables demográficas
- Clasificar imágenes satelitales
- Identificar tópicos en corpus textuales
- Clasificar textos según su sentimiento -positivos o negativos-
- Detectar objetos en imágenes

# ¿Qué es Machine Learning?

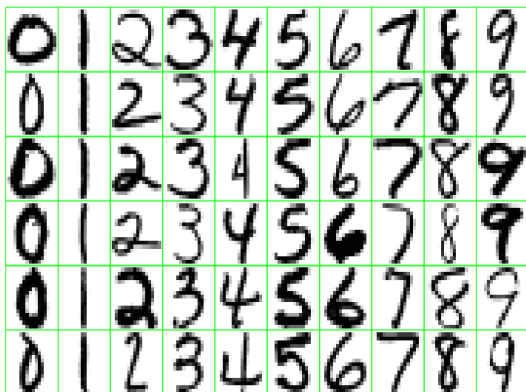
## Ataques cardíacos



- Hastie, et al. 2013 [8]

# ¿Qué es Machine Learning?

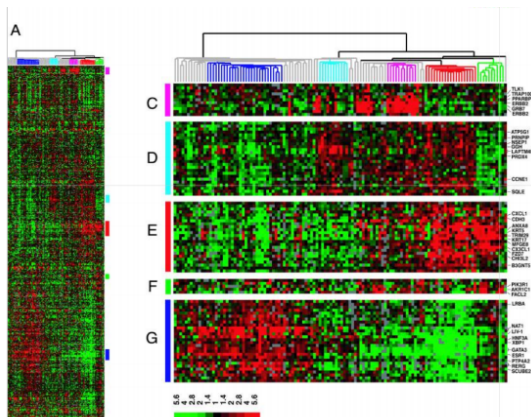
## Reconocimiento de dígitos



- Hastie, et al. 2013 [8]

# ¿Qué es Machine Learning?

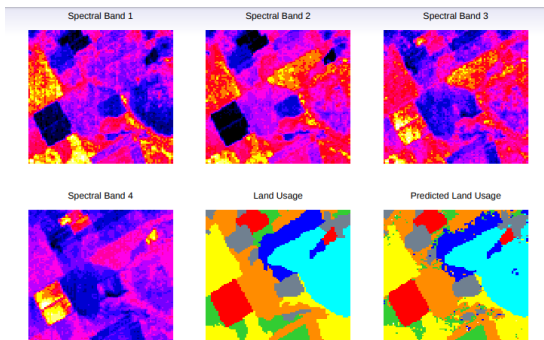
## Cáncer y genes



- Hastie, et al. 2013 [8]

# ¿Qué es Machine Learning?

## Imágenes satelitales

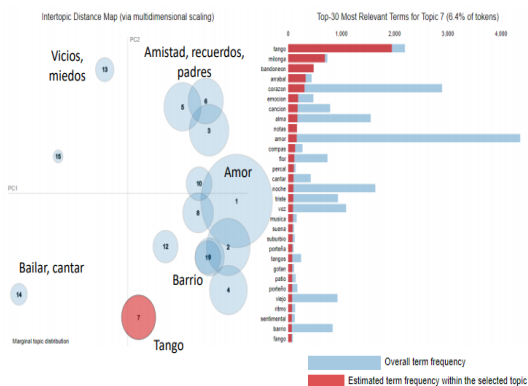


*Usage  $\in$  {red soil, cotton, vegetation stubble, mixture, gray soil, damp gray soil}*

- Hastie, et al. 2013 [8]

# ¿Qué es Machine Learning?

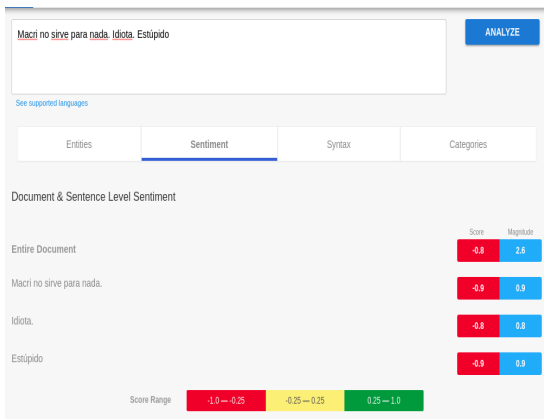
## Detección de tópicos -tango-





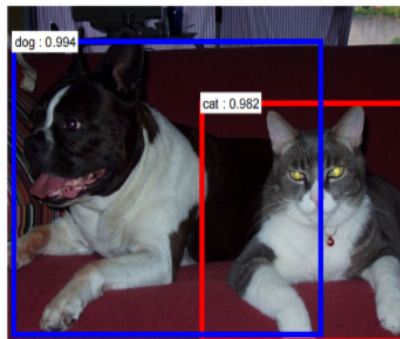
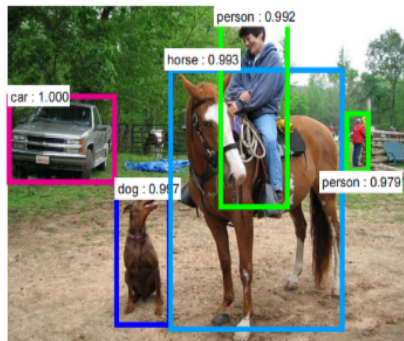
# ¿Qué es Machine Learning?

## Sentiment Analysis



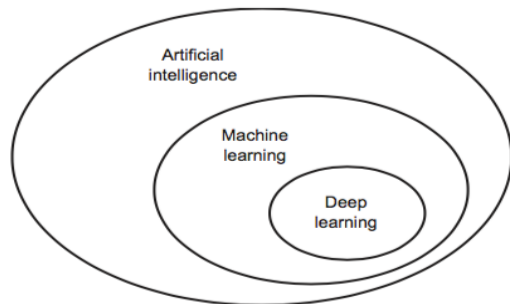
# ¿Qué es Machine Learning?

## Detección de objetos en imágenes



# Delimitación del campo

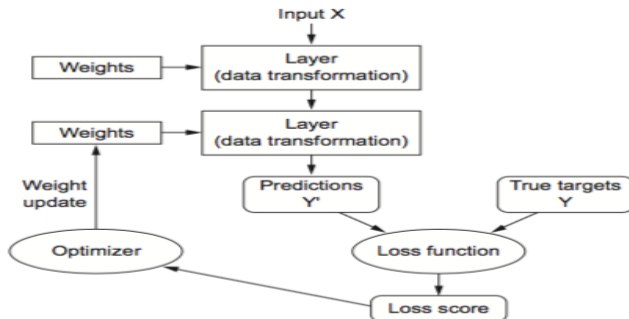
Chollet, 2017 [4]



**Figure 1.1** Artificial intelligence, machine learning, and deep learning

# Delimitación del campo

Chollet, 2017, [4]



## Aprendizaje Supervisado:

- Variable dependiente, resultado, target  $Y$
- Matriz de  $p$  predictores  $X$ , featuras, variables independientes, etc...
- En *problemas de regresión*  $Y$  es cuantitativa
- En *problemas de clasificación*  $Y$  es cualitativa
- Tenemos datos de entremiento  $(X_1, y_1, \dots, (X_N, y_1)$ . Son observaciones (ejemplos, instancias) de las mediciones

## Aprendizaje No Supervisado:

- No hay una variable target... no existe  $Y$
- Sí existe la matriz de  $p$  predictores  $X$ , features, variables independientes, etc...
- Es más complejo: es más difícil medir qué tan bien funciona el modelo

# ¿Qué es un modelo?

- Básicamente: una manera de proponer hipótesis sobre la forma en que se combinan variables
- En general, vamos a estar tratando de generar modelos de esta forma

$$Y = f(X) + \epsilon \quad (1)$$

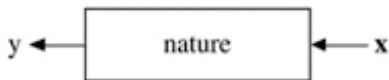
- Todo el problema es estimar  $f(X)$ , es decir, de qué forma(s) se combinan las  $X$  para generar un output
- Una posibilidad es suponer que  $Y$  es una combinación lineal de las  $X$

# ¿Qué es un modelo?

Las dos culturas (Breiman, 2001) [3]

*“Todos los modelos son equivocados. Algunos son útiles.”*

*George Box*



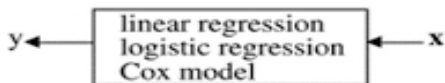
- El mundo como productor de outputs  $-y-$  en base a features  $-X-$
- Problemas: ¿cuál es la manera en que el mundo produce resultados?
- Una forma común es asumir que los datos son generados por extracciones independientes de  
 $output = f(\text{predictores}, \text{ruido}, \text{parametros})$



# ¿Qué es un modelo?

Las dos culturas (Breiman, 2001)[3]

## Modelado estadístico

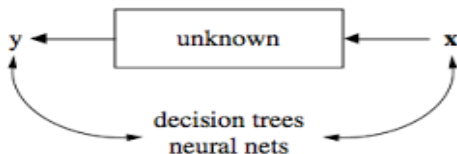


- Énfasis en  $f(x)$ . El modelo se postula en base a supuestos sobre  $f(x)$
- Conocimiento acumulado, teoría, diseño de experimentos
- Los parámetros son estimados con los datos y luego se realizan las predicciones.
- Evaluación del modelo: estimadores insesgados, robustos, mínima varianza

# ¿Qué es un modelo?

Las dos culturas (Breiman, 2001)[3]

## Modelado algorítmico (o *Machine Learning*, *Data Mining*, etc.)



- Énfasis en  $\hat{y}$
- El enfoque es encontrar una función  $f(x)$  -un algoritmo- que opera sobre las  $x$  para predecir las  $y$ .
- El modelo se “aprende” de los datos
- Evaluación del modelo: performance predictiva

# ¿Qué es un modelo?

El fin de la teoría (Anderson, 2008)[1]

CHRIS ANDERSON SCIENCE 06.23.08 12:00 PM

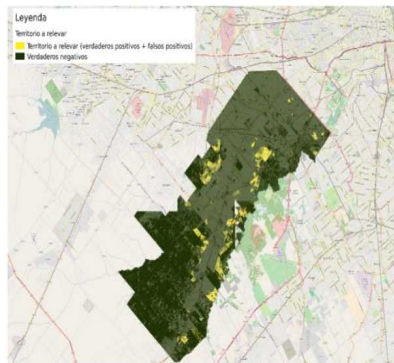
## THE END OF THEORY: THE DATA DELUGE MAKES THE SCIENTIFIC METHOD OBSOLETE



*Out with every theory of human behavior, from linguistics to sociology. Forget taxonomy, ontology, and psychology. Who knows why people do what they do? The point is they do it, and we can track and measure it with unprecedented fidelity. With enough data, the numbers speak for themselves.*

# Aplicaciones de ML en Ciencias Sociales

## Detección de asentamientos informales en base a imágenes satelitales



- Baylé, 2016 [2]

# Aplicaciones de ML en Ciencias Sociales

## Imputación de datos perdidos en encuestas sociodemográficas

- Rosati, 2017 [10]

ISSN 1813-4121

SaberEs

REVISTA DE CIENCIAS ECONÓMICAS Y ESTADÍSTICAS

INICIO ACERCA DE INDICAR SESIÓN REGISTRARSE BUSCAR ACTUAL ARCHIVOS ARTÍCULOS

Inicio > Vol. 9, Núm. 1 (2017) > Rosati

Construcción de un modelo de imputación para variables de ingreso con valores perdidos a partir de ensemble learning. Aplicación en la Encuesta Permanente de Hogares (EPH)

Germán Federico Rosati

Resumen

El presente documento se propone exponer los avances realizados en la construcción de un modelo de imputación de valores perdidos y sin respuesta para las variables de ingreso en encuestas a hogares. Se presentará la propuesta metodológica general y los resultados de las pruebas realizadas. Se evalúan dos tipos de modelos de imputación de datos perdidos: 1) el método hot deck (implementado utilizado por organismos importantes en el Sistema Estadístico Nacional, tales como la Encuesta Permanente de Hogares y la Encuesta Anual de Hogares de la Ciudad de Buenos Aires) y 2) un ensemble de modelos de regresión LASSO (Least Absolute Shrinkage and Selection Operator). El mismo se basa en la generación de múltiples modelos de regresión LASSO a través del algoritmo bagging y de su agregación para la generación de la imputación final. En la primera y segunda parte del documento plantea el problema de forma más específica y se pasa revista a los principales mecanismos de generación de los valores perdidos y sus implicaciones que los mismos tienen al momento de generar modelos de imputación. En el tercer apartado se describen los métodos de imputación más habitualmente utilizados, enfatizando sus ventajas y limitaciones. En la cuarta parte, se desarrollan los fundamentos teóricos y metodológicos de los dos técnicas de imputación propuestas. Finalmente, en la quinta sección, se presentan algunos resultados de la aplicación de los métodos propuestos a datos de la Encuesta Permanente de Hogares.

SciELO  
Revista de Ciencias Económicas y Estadísticas

latindex  
REVISTA DE CIENCIAS ECONÓMICAS Y ESTADÍSTICAS

NB  
Núcleo de Bibliotecas

REDIB  
Red Iberoamericana

Dialnet

# Aplicaciones de ML en Ciencias Sociales

## Integración de comunidades migrantes

- Lamanna, Lenormand, et al 2016 [9]

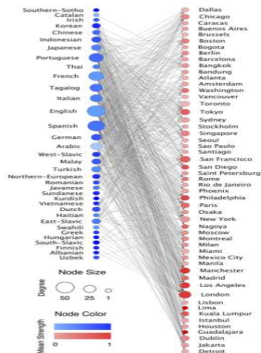
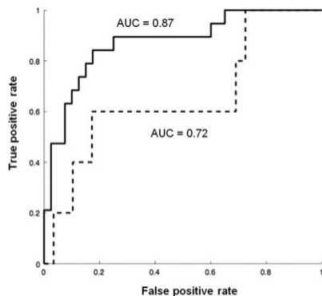


Fig 2. Bipartite spatial integration net work. The network comprises of two sets: L of Languages and C of cities; the languages detected are connected to the cities set where the corresponding community of immigrants has been found. The weight of the edge corresponds to the values of  $A_{ij}$ . The size of the nodes is proportional to its degree and the color to its mean strength.

# Aplicaciones de ML en Ciencias Sociales

## Predicción de enfermedades mentales mediante text mining



**Figure 2** Receiver operating characteristics (ROC) for the University of California Los Angeles (UCLA) clinical high-risk (CHR) classifier of psychosis outcome as applied to the UCLA dataset (solid line) and to the realigned New York City (NYC) dataset (dotted line). AUC – area under the curve.

- Corcoran, Carrillo, Fernández Slezak et al, 2018 [5]

# Aplicaciones de ML en Ciencias Sociales

## Automatización de procesos para la construcción de bases de datos de protestas

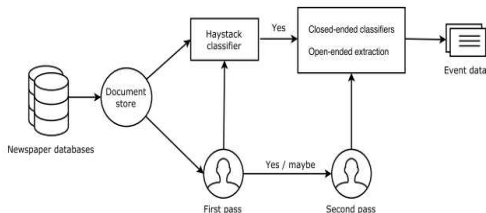


Figure 1: MPEDS pipeline with training.

- Hannah, 2017 [7]



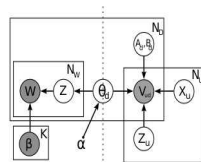
# Aplicaciones de ML en Ciencias Sociales

## Posiciones ideológicas en proyectos de ley

- Gerrish y Blei, 2015 [6]

Terrorism	Commemorations	Transportation
terrorist	nation	transportation
september	people	minor
attack	life	print
nation	world	tax
york	serve	land
terrorist attack	percent	guard
hezbollah	community	coast guard
national guard	family	substitute

Labeled topics

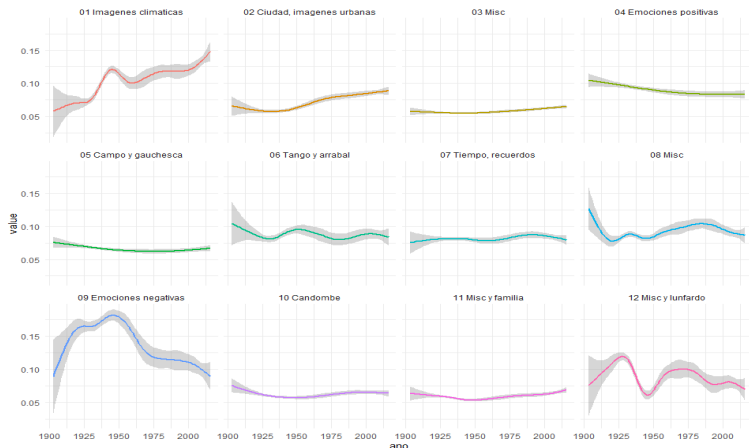


The issue-adjusted ideal point model

Figure 3: Left: Top words from topics fit using labeled LDA [6]. Right: the issue-adjusted ideal point model, which models votes  $v_{ud}$  from lawmakers and legislative items. Classic item response theory models votes  $v$  using  $x_u$  and  $a_d, b_d$ . For our work, documents' issue vectors  $\theta$  were estimated fit with a topic model (left of dashed line) using bills' words  $w$  and labeled topics  $\beta$ . Expected issue vectors  $\mathbb{E}_q[\theta|w]$  are then treated as constants in the issue model (right of dashed line).

# Aplicaciones de ML en Ciencias Sociales

Tópicos en letras de tango, Rosati, 2020, inédito

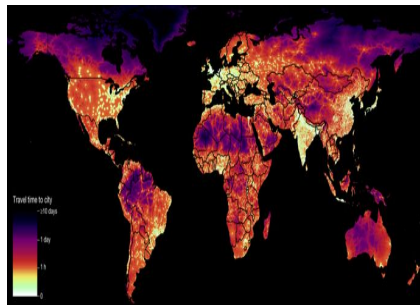


# Pero antes..

Algunos ejemplos de textos con problemas relevantes

## Datos de accesibilidad como acercamiento a la desigualdad

Weiss et al 2018 [12]





# Referencias bibliográficas I



ANDERSON, C.

The End of Theory: The Data Deluge Makes the Scientific Method Obsolete.  
*Wired* 16, 07 (June 2008).



BAYLÉ, F.

Detección de villas y asentamientos informales en el partido de la matanza mediante teledetección y sistemas de información geográfica.  
Master's thesis, Tesis de Maestría, 2016.



BREIMAN, L.

Statistical modeling: The two cultures (with comments and a rejoinder by the author).  
*Statistical Science* 16, 3 (08 2001), 199–231.



CHOLLET, F.

*Deep Learning with Python*.  
Manning, Nov. 2017.



CORCORAN, C. M., CARRILLO, F., FERNÁNDEZ-SLEZAK, D., BEDI, G., KLIM, C., JAVITT, D., BEARDEN, C., AND CECCHI, G.

Prediction of psychosis across protocols and risk cohorts using automated language analysis.

*World Psychiatry* 17, 01 (2 2018).

# Referencias bibliográficas II



GERRISH, S., AND BLEI, D. M.

How they vote: Issue-adjusted models of legislative behavior.

In *NIPS* (2012), P. L. Bartlett, F. C. N. Pereira, C. J. C. Burges, L. Bottou, and K. Q. Weinberger, Eds., pp. 2762–2770.



HANNA, A.

Mpeds: Automating the generation of protest event data, Jan 2017.



HASTIE, T., TIBSHIRANI, R., AND FRIEDMAN, J.

*The Elements of Statistical Learning*.

Springer Series in Statistics. Springer New York Inc., New York, NY, USA, 2001.



LAMANNA, F., LENORMAND, M., SALAS-OLMEDO, M. H., ROMANILLOS, G., GONÇALVES, B., AND RAMASCO, J. J.

Immigrant community integration in world cities.

*CoRR abs/1611.01056* (2016).



ROSATI, G.

Construcción de un modelo de imputación para variables de ingreso con valores perdidos a partir de Ensamble Learning. Aplicación a la Encuesta Permanente de Hogares (EPH).

*SaberES. Revista de Ciencias Económicas y Estadística* 9, 01 (Febrero 2017).

# Referencias bibliográficas III



SANZ, C., ZAMBERLAN, F., EROWID, E., EROWID, F., AND TAGLIAZUCCHI, E.

The experience elicited by hallucinogens presents the highest similarity to dreaming within a large database of psychoactive substance reports.

*Frontiers in Neuroscience* 12 (2018), 7.



WEISS, D. J., NELSON, A., GIBSON, H. S., TEMPERLEY, W., PEEDELL, S., LIEBER, A., HANCHER, M., POYART, E., BELCHIOR, S., FULLMAN, N., MAPPIN, B., DALRYMPLE, U., ROZIER, J., LUCAS, T. C. D., HOWES, R. E., TUSTING, L. S. AND KANG, S. Y., CAMERON, BISANZIO, D., BATTLE, K. E., BHATT, S., AND GETHING, P. W.

A global map of travel time to cities to assess inequalities in accessibility in 2015.

*Nature* 553, 333 (01 2018).