

# Unidad: Introducción a la Ciencia de Datos

## Ciencia de Datos para Economía y Negocios

---

Nicolás Sidicaro

Agosto 2025

# Nicolás Sidicaro

- Investigador en Fundar
- Docente FCE-UBA y UADE
- Licenciado en Economía
- Data Scientist
- Econometría y minería de datos

## Contacto

- Email: nsidicaro.fce@gmail.com

- Hacia un dispositivo público de monitoreo del subte de Buenos Aires
  - Monitor de fallas del subte mediante API de Twitter
- Distribución de trabajadores en el AMBA
  - Georreferenciación
- Mapa productivo laboral de Argentina
  - Georreferenciación y visualización
- El precio de la ropa en Argentina
  - Relevamiento de precios

# 1. Ciencia de Datos: Introducción

## ¿Qué es la Ciencia de Datos?

- Campo interdisciplinario que utiliza métodos científicos, procesos, algoritmos y sistemas para extraer conocimiento e información de datos estructurados y no estructurados
- Combina aspectos de:
  - Estadística
  - Matemáticas
  - Programación
  - Visualización
  - Conocimiento del dominio
- Se enfoca en descubrir patrones, tendencias y relaciones para generar información accionable

# 1. Ciencia de Datos: Introducción

## ¿Para qué se puede usar?

- Toma de decisiones basada en datos
- Optimización de procesos y recursos
- Detección de anomalías y fraudes
- Personalización de productos y servicios
- Predicción de tendencias y comportamientos
- Automatización de procesos

# Ecosistema de Ciencia de Datos

## Data Mining (Minería de Datos)

- Proceso de descubrir patrones y conocimientos interesantes a partir de grandes volúmenes de datos
- Se enfoca en la **extracción** de información no evidente
- Utiliza algoritmos para identificar relaciones, anomalías y tendencias

## Data Analysis (Análisis de Datos)

- Proceso de inspección, limpieza, transformación y modelado de datos
- Se enfoca en el **examen** de datos para responder preguntas específicas
- Más orientado a la comprensión descriptiva y diagnóstica

## Arquitectura de Datos

- Diseño de estructuras para recopilar, almacenar, procesar y consumir datos
- Definición de flujos de datos, bases de datos y sistemas de procesamiento
- Garantiza que los datos estén disponibles, seguros y sean de calidad

# Ecosistema de Ciencia de Datos

## Machine Learning (Aprendizaje Automático)

- Subcampo de la inteligencia artificial que permite a los sistemas aprender de datos
- Crea modelos para reconocer patrones y tomar decisiones con mínima intervención humana
- **Tipos principales:**
  - Supervisado: aprende de datos etiquetados
  - No supervisado: encuentra patrones en datos no etiquetados

## Inteligencia Artificial (IA)

- Campo más amplio que busca crear sistemas que puedan percibir, razonar y actuar
- Machine Learning es un subconjunto de la IA
- Incluye procesamiento de lenguaje natural, visión por computadora, robótica, etc.
- Enfocada en crear soluciones que emulan aspectos de la inteligencia humana

# Comparación de conceptos clave

Concepto	Definición	Enfoque	Aplicación en economía	Uso en negocios
Ciencia de datos	Campo interdisciplinario que extrae conocimiento y valor de los datos	Proceso completo desde preguntas hasta decisiones	Análisis integral de problemas económicos complejos	Transformación digital, innovación basada en datos, optimización de procesos empresariales
Análisis de datos	Examen de datos para sacar conclusiones sobre la información	Enfoque principalmente descriptivo y diagnóstico	Interpretación de indicadores económicos y tendencias	Informes de desempeño, dashboards de KPIs, análisis de competencia

# Comparación de conceptos clave

Concepto	Definición	Enfoque	Aplicación en economía	Uso en negocios
Data Mining	Descubrimiento de patrones en grandes conjuntos de datos	Exploratorio, busca relaciones no evidentes	Segmentación de mercados, patrones de consumo	Análisis de canasta de mercado, sistemas de recomendación, segmentación de clientes
Machine Learning	Algoritmos que mejoran automáticamente con la experiencia	Predictivo y prescriptivo	Predicción de variables económicas, detección de anomalías	Predicción de demanda, detección de fraude, mantenimiento predictivo



# Comparación de conceptos clave

Concepto	Definición	Enfoque	Aplicación en economía	Uso en negocios
Inteligencia Artificial	Sistemas que emulan comportamiento inteligente	Resolución autónoma de problemas complejos	Automatización de decisiones económicas complejas	Asistentes virtuales, automatización de servicio al cliente, optimización logística

# Etapas del proceso de Ciencia de Datos

1. **Definición del problema:** Identificar objetivos y preguntas clave
2. **Recolección de datos:** Obtener información relevante de diversas fuentes
3. **Limpieza y preparación:** Transformar datos crudos en formato utilizable
4. **Exploración y análisis:** Identificar patrones y relaciones
5. **Modelado:** Crear modelos predictivos o descriptivos
6. **Evaluación e interpretación:** Validar resultados y extraer conclusiones
7. **Implementación y comunicación:** Aplicar hallazgos y presentarlos efectivamente

El 80% del tiempo de trabajo suele estar en las primeras 3 etapas. La parte **fancy** es marginal.

# 2. Herramientas fundamentales: GitHub

## ¿Qué es GitHub?

- Plataforma basada en Git para control de versiones y colaboración
- Permite a múltiples personas trabajar en los mismos archivos sin conflictos
- Funciona como un "repositorio" central para código y documentación

## Funcionalidades principales

- **Repositorios:** Almacenamiento de proyectos con historial completo
- **Branches (Ramas):** Versiones paralelas para desarrollo simultáneo
- **Pull Requests:** Mecanismo para revisar y aprobar cambios
- **Issues:** Seguimiento de tareas, errores y funcionalidades
- **Actions:** Automatización de flujos de trabajo

# 2. Herramientas fundamentales

## ¿Qué es Google Colaboratory (Colab)?

- Entorno de notebook basado en Jupyter en la nube
- Permite escribir y ejecutar código Python/R directamente en el navegador
- No requiere instalación ni configuración local

## Características principales

- **Integración con Google Drive:** Almacenamiento y acceso a datos
- **GPU/TPU gratuitas:** Aceleración para modelos complejos
- **Entorno preconfigurado:** Bibliotecas populares ya instaladas
- **Interfaz interactiva:** Combina código, texto narrativo y visualizaciones
- **Fácil compartición:** Colaboración en tiempo real

# 2. Herramientas fundamentales: R

## ¿Qué es R?

- Lenguaje de programación especializado en computación estadística y gráficos
- Software libre y de código abierto
- Creado por Ross Ihaka y Robert Gentleman en 1993
- Ampliamente utilizado en investigación estadística, ciencia de datos y machine learning

## Características principales

- **Orientado al análisis estadístico:** Diseñado específicamente para esta tarea
- **Extensible:** Más de 18,000 paquetes adicionales en CRAN
- **Capacidades gráficas avanzadas:** Excelente para visualización de datos
- **Comunidad activa:** Amplio soporte y recursos disponibles
- **Reproducibilidad:** Facilita documentar y compartir análisis completos

## Paquetes esenciales que utilizaremos

tidyverse: dplyr; ggplot2; tidyr; plotly, data.table, lubridate, caret; rvest, RSelenium, haven

# 2. Herramientas fundamentales: RStudio

## ¿Qué es RStudio?

- Entorno de desarrollo integrado (IDE) para R
- Interfaz gráfica que facilita el uso de R
- Disponible en versión de escritorio y servidor
- Desarrollado por Posit (anteriormente RStudio, Inc.)

## Componentes principales

- **Editor de código:** Con resaltado de sintaxis y autocompletado
- **Consola:** Para ejecutar comandos de R
- **Entorno y variables:** Visualización de datos y objetos en memoria
- **Historial:** Registro de comandos ejecutados
- **Gráficos:** Visualización de resultados
- **Ayuda y documentación:** Acceso rápido a información sobre funciones

# 2. Herramientas fundamentales: Python

## ¿Qué es Python?

- Lenguaje de programación de alto nivel, interpretado y de propósito general
- Diseño centrado en la legibilidad del código
- Multiparadigma: soporta programación orientada a objetos, imperativa y funcional
- Uno de los lenguajes más populares para ciencia de datos y machine learning

## Características principales para ciencia de datos

- **Sintaxis clara y legible:** Facilita el aprendizaje y mantenimiento
- **Ecosistema científico robusto:** NumPy, pandas, SciPy, scikit-learn, etc.
- **Visualización potente:** Matplotlib, Seaborn, Plotly
- **Versatilidad:** Se integra fácilmente con otros lenguajes y sistemas
- **Machine learning y deep learning:** Bibliotecas como scikit-learn, TensorFlow, PyTorch

## Aplicaciones en nuestro curso

- Implementación de algoritmos de machine learning

## 2. ¿Cuál uso?

Tipo de preguntas del tipo ¿Dios existe? ¿Messi o Maradona? No hay una respuesta clara

### Ventajas R:

- Gráficos más lindos
- Un poco más rápido
- Más enfocado en estadística

### Ventajas Python

- Mejor integración para otras cuestiones de programación (se puede usar Python para otras cosas además del análisis de datos)
- Más robusto en Machine Learning

**Entonces, hay que usar los dos** y quizás un poco STATA para algunas cosas de econometría



# 3. Estructura del curso

## Modalidad de dictado

- **Clases presenciales** (martes): enfoque conceptual y teórico
- **Clases virtuales** (viernes): enfoque práctico
  - Algunas sincrónicas, otras asincrónicas según cronograma
- Complemento con lecturas de bibliografía y ejercicios prácticos

## Tres grandes unidades

1. **Programación** orientada a análisis de datos
2. **Estadística** aplicada al contexto
3. **Visualización** de datos y resultados

# 3. Estructura del curso: Evaluación

**Primer parcial:** se evaluarán capacidades de entendimiento de los procesos lógicos de programación (pseudocódigo) - *Presencial (multiple choice individual)*

**Segundo parcial:** se evaluarán los conocimientos de estadística aplicados a ejercicios - *Domiciliario (individual)*

**Tercer parcial:** se evaluará el uso de visualizaciones para una base de datos. Formato competencia/hackaton - *Domiciliario (grupal)*

**Trabajo final:** análisis de una base de datos a elección. Las bases disponibles para seleccionar se darán a conocer luego del primer parcial.

**Recuperatorio:** se puede recuperar el primer parcial o el segundo (UNO DE LOS DOS) de forma presencial el 18/11/2025

# 3. Estructura del curso: Evaluación

**Aprobación del curso:** se deben aprobar todos los parciales y el trabajo final para aprobar el curso. La ponderación es 50% del Trabajo Final y 50% de los parciales. La nota de parciales es una ponderación de cada uno con los siguientes pesos: 40% el primero; 30% el segundo y 30% el tercero

**¿Qué pasa si no puedo entregar el trabajo final esa semana?:** hay una semana más para entregarlo, pero cada día que pasa desde la entrega original equivale a 0,3 puntos menos (hasta que se llega a dos puntos menos)

**No hay final**

# 3. Estructura del curso: Los Grupos

- Inicialmente está todo pensado para que los grupos sean de a dos
- Sin embargo, siempre hay quien no lo quiere hacer de a dos...
- Para estas navidades llega: la **curva de la complejidad**

# 3. Estructura del curso: Los Grupos

- Grupo individual: corrección muy rigurosa por parte del docente
- Grupo de a dos: corrección amable (original)
- Grupo de a tres: corrección un poco más rigurosa, pero no tanto
- Grupo de a cuatro: corrección muy rigurosa por parte del docente

# 3. Estructura del curso: Los Grupos

- Se van a tener que anotar en un Google Form
- No se van a poder repetir las bases de datos más de 2 grupos, por lo que van a tener que revisar un archivo que les voy a compartir en el momento adecuado y elegir en función de eso. Si eligen una base que ya estaba completa, les voy a avisar y van a tener que elegir otra.

# 3. Estructura del curso: Materiales

## Herramientas principales

- **R o Python** indistintamente. Las clases las voy a dar en R, pero va a haber material complementario que trata esos mismos temas en Python.
- Los **parciales no requieren un lenguaje en particular**, así que si usan otros (STATA, Julia, algún otro) no hay problema.
- Si hacen algún código y quieren que lo vea o no les sale algo, podemos verlo en horario de clase. También voy a habilitar la próxima semana un canal para que carguen sus códigos y voy a revisar al azar (opcional)

# 4. Recursos

- **Material de trabajo:**
  - Slides
  - Scripts compartidos en GitHub (desde clase 2)
  - Datasets de práctica en Github (desde clase 2)
  - Recursos online (tutoriales, documentación)
  - Bibliografía complementaria



# 4. Aplicaciones en Economía y Negocios

## Economía

- Análisis de tendencias macroeconómicas
- Predicción de indicadores económicos
- Evaluación de impacto de políticas públicas
- Estudios de comportamiento del consumidor
- Análisis de mercados laborales (con EPH)
- Comercio internacional y análisis de exportaciones

## Negocios

- Segmentación de clientes
- Análisis predictivo de ventas
- Optimización de precios
- Detección de fraude
- Análisis de sentimiento en redes sociales
- Optimización de cadenas de suministro
- Sistemas de recomendación

# ¡Gracias!

Contacto: [nsidicaro.fce@gmail.com](mailto:nsidicaro.fce@gmail.com)