SPRINT 1 - ENCUENTRO 1 FEB Introducción a Data Science -**Python Dashboard** 

**Sprint 1** 

?

### Hacia la utilización de los datos para mejorar la toma de decisiones "The world is one big data problem." - Andrew McAfee, co-director of the MIT

El equipo docente subió contenido extra del encuentro 🛠

**☐** Grabación

Presentación

Lista de recursos

Glosario

Profundiza

Notebook:

Redes

Contenido extra

DS\_Bitácora\_01\_Ejemplo

Potencia tu Talento -

DS-ONLINE-6 2/2/21

for Digital Business at the MIT Sloan School of Management.

Initiative on the Digital Economy, and the associate director of the Center

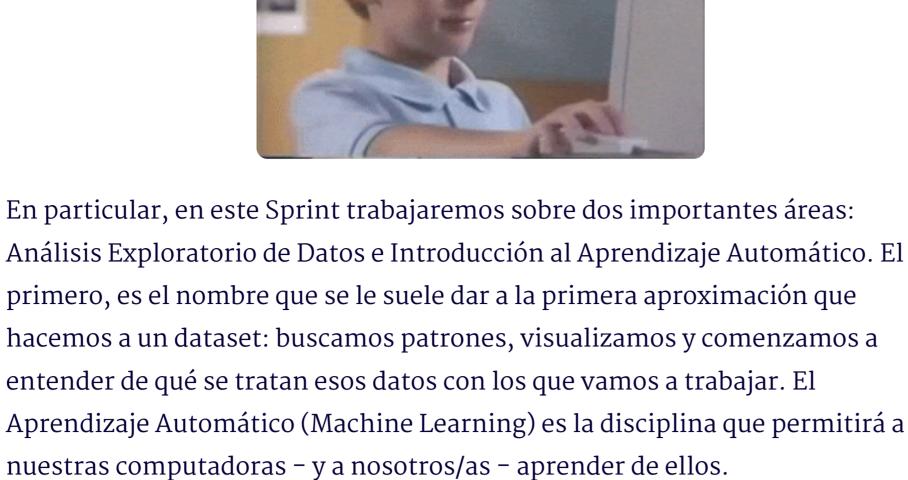
## Los primeros pasos en Data Science

¡Bienvenido/a a tu primera bitácora! Aquí encontrarás disponibles los

tu acceso a los canales de comunidad: te propondremos interactuar con otros/as Data Scientists para disipar dudas, hacer consultas y, por qué no, conocer y recomendar cosas interesantes sobre los temas que trabajaremos. El primer Sprint de Aprendizaje (que consta de 16 bitácoras para 16 encuentros), tendrá por objetivo familiarizarnos con las herramientas fundamentales de Data Science, lo que nos permitirá comenzar a pensar como Data Scientists.

conceptos clave y las herramientas que serán necesarias para prepararte para

cada encuentro con tu equipo docente y compañeros/as. Además, motivaremos



Para lograr estos objetivos, estudiaremos: • Un lenguaje de programación (Python) junto con las librerías propias del análisis de datos. • Probabilidad y Estadística: repasaremos y estudiaremos las herramientas de estas disciplinas que son fundamentales para el análisis de datos

**ANALIZAR** 

**PRESENTAR** 

- ¿Qué es Data Science (o Ciencia de Datos)?
- Nuestro tiempo de ocio, que muchas veces es poco pero no por eso menospreciado, queremos aprovecharlo al máximo. Sentarnos en el sillón a

mirar una película o una serie en Netflix puede que sea todo lo que necesitamos

un sábado a la noche después de tanta actividad. Pero a veces nos pasa que la

perdernos dentro de los cientos de títulos de los cuales podemos elegir. Entrar

en la paradoja de la elección puede hacernos perder mucho tiempo, y hasta

siempre. Pero, por suerte, (jo gracias a los/as Data Scientists!) Netflix y otras

plataformas de streaming tienen opciones para recomendarnos en base a

plataforma cuenta con tantas opciones que lentamente comenzamos a

sacarnos las ganas de ver algo nuevo para terminar viendo lo mismo de

#### nuestra interacción con ellas. ¿Cómo lo hacen? Netflix cuenta con una gran cantidad de información de sus usuarios: qué vieron, qué les gustó y qué no, qué cosas dejaron sin terminar, etc. Con esa información, puede decidir qué

recomendarnos en base a nuestro historial como usuarios, pero también a partir de otros/as usuarios/as que vieron contenido igual o similar al que vimos nosotros. Cómo hacer esas recomendaciones es un problema de Data Science. Data Science es un campo **interdisciplinario** tanto en sus objetivos como en sus metodologías que, por medio de un proceso iterativo, busca: **DEFINIR** Partimos de preguntas que queremos responder. ¿Cuáles son los datos que necesitamos para responder las preguntas que nos hacemos? **INVESTIGAR** Se obtienen los datos, se "limpian" y se procede a explorarlos.

¿Qué hace un/a Data Scientist?

Esta podría ser una primera aproximación a la respuesta: un/a data scientist

• Dado un problema, ¿qué datos nos pueden ayudar a abordarlo y resolverlo?

(científico/a de datos) es alguien que, se hace tres preguntas:

• ¿Dónde – y cómo – podemos conseguir datos?

Presentamos los resultados obtenidos y las conclusiones a las que llegamos.

Los datos obtenidos se analizan con modelos (estadísticos, Machine Learning,

etc.). Interpretamos los resultados y transformamos datos en información.

#### trabajar con él? En este sentido, el diagrama de Venn de Drew Conway ubica a la tarea de un

Data Scientist en la intersección entre los siguientes conocimientos y

• Dado un conjunto de datos, ¿qué problemas interesantes podemos

competencias:

- Conocimiento **Especializado**

• Podríamos agregar las habilidades para comunicar los resultados.

La descripción que aporta Favio Vázquez resulta esclarecedora:

para crear modelos y responder la pregunta inicial".

trabajar con datos?

¿Qué es programar?

"un Científico de Datos es una persona (...) encargada de analizar problemas de negocios/organizaciones y ofrecer una solución estructurada que comienza convirtiendo este problema en una pregunta válida y completa. Luego, utilizando el método científico, la programación y las herramientas

computacionales desarrollar códigos que preparan, limpian y analizan datos

¿Por qué programar si elegimos

encontraste en un cuarto oscuro de revelado. La extracción de datos consiste en utilizar los equipos para revelar las fotos lo más rápido posible, para que puedas ver si hay algo inspirador o interesante en ellas. Al igual que con las fotos, recuerda no tomarte en serio lo que ves. Tú no tomaste las fotos, así que no sabes mucho sobre las historias que hay detrás de ellas. La regla de oro de la minería de datos es: enfocarse en lo que está aquí" (Cassie Kozyrkov — Head of Decision Intelligence, Google).

Siguiendo el ejemplo de Cassie Kozyrkov, los lenguajes de programación -como

Python o R-, serían la herramienta de revelado. ¡Con calma! Si bien vamos a

usar un lenguaje y aprenderemos a programar, esta NO es una carrera de

"Piensa en tu conjunto de datos como un grupo de fotos en negativo que las

#### programación. Aprenderemos lo necesario para resolver lo que un/a data scientist necesita. Los grandes volúmenes de datos que existen actualmente y los métodos utilizados para analizarlos nos exigen trabajar de manera eficiente pero flexible. Esto sólo es posible hacerlo utilizando un lenguaje de programación.

Cuando buscamos definiciones acerca de qué es la programación, nos

encontramos con respuestas tales como "dar instrucciones al ordenador" o

bien "crear software usando un lenguaje de programación". Según wikipedia,

es el proceso por el cual una persona desarrolla un programa valiéndose de una herramienta que le permita escribir el código (nosotros usaremos Python) y de otra que sea capaz de "traducirlo" a lo que se conoce como lenguaje de máquina, que puede comprender el ordenador. ¿Por qué Python? Python es un lenguaje de programación de código abierto (open source) con una

sintaxis simplificada y flexible que lo convierte en una de las mejores opciones

• Es fácil de usar, sobre todo cuando lo comparamos con otros lenguajes de

para utilizar. Existen muchos motivos para elegir Python:

**Primeros Pasos con Python** 

#### • Hay una gran comunidad usándolo, en particular en Data Science, Aprendizaje Automático e Inteligencia Artificial. • Cuenta con una amplia cantidad de **librerías** específicas (¡pronto veremos

programación.

qué son!).

• Es rápido y eficiente.

jerga que puedes encontrar es:

lenguaje de programación interpretado).

también, de la tarea a desarrollar.

los/as usuarios/as y programadores/as, mientras que los de bajo nivel están más próximos a la arquitectura del hardware. En general, existe un compromiso entre Alto y Bajo nivel y la velocidad de ejecución, su eficiencia, etc.

• Compilados o Interpretados: cuando programamos, escribimos

escritas en un lenguaje que el humano comprende, deben ser

instrucciones que la computadora deberá ejecutar. Estas instrucciones,

computadora entienda (en el fondo, o y 1). Un lenguaje compilado requiere

convertir el código a lenguaje de máquina. Un lenguaje interpretado, en

cambio, es convertido a lenguaje de máquina a medida que es ejecutado.

Python es un lenguaje de **alto nivel** (¡aunque no lo creas!) e **interpretado**. En

general, en Data Science se utilizan lenguajes interpretados (R también es un

Sigamos avanzando. Si programar es "dar instrucciones al ordenador", de

alguna forma debemos dar esas instrucciones. En la programación tradicional,

esas instrucciones están escritas en uno o varios archivos (scripts), que juntos

"traducidas" al lenguaje de máquina, es decir, instrucciones que la

un paso adicional antes de ser ejecutado —la compilación— para

• Bajo o Alto nivel: los lenguajes de alto nivel son más comprensibles para

Hay muchas formas de agrupar a los lenguajes de programación. Un poco de la

forman un programa. Para escribir esas instrucciones, el código es escrito usando un entorno de desarrollo, que es simplemente la interfaz con la computadora que usamos para, precisamente, escribir código. Para darte una idea, un entorno de desarrollo puede ser el Bloc de Notas de Windows, aunque en general usamos entornos más avanzados con muchas funcionalidades (colores específicos para instrucciones o palabras clave, autocompletar, compilar o ejecutar el código a medida que lo vamos desarrollando, etc.). Qué

entorno de desarrollo utilizaremos dependerá del lenguaje de programación y,

En Data Science el entregable no suele ser un programa, pero tampoco un

Notebook: DS\_Bitácora\_01\_Ejemplo

Con el fin de generar estos notebooks, contamos con varios entornos de

desarrollo. El más conocido es **Jupyter Notebook**, que instalaremos en

nuestras computadoras. Cuando deseemos trabajar en la nube, usaremos

Google Colab (que también corre sobre Jupyter Notebook). Sigue las

informe. Sino, más bien, lo que se llama un **notebook**. Un notebook es la combinación de código con informe. Incluye el código utilizado para analizar los datos o generar figuras, intercalado con texto en lenguaje natural mediante el que explicamos qué estamos haciendo, los resultados alcanzados, las conclusiones y hasta ecuaciones matemáticas. Aquí va un ejemplo:

# Instala Jupyter

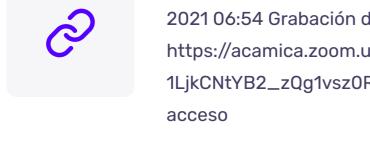
¡Prepárate para el próximo encuentro!

### Glosario Palabras claves para convertirte en un experto.

Potencia tu Talento - Redes Tips para construir tu perfil profesional.

# **DS-ONLINE-6 2/2/21**

Contenido extra del encuentro



instrucciones que te dejamos en la sección de "Herramientas". ¡Lo necesitarás para poner manos a la obra durante el encuentro! Advertencia: hay una altísima probabilidad de que te sientas confundido/a con

la cantidad de nombres nuevos, con la jerga, etc. No te preocupes, jes común! Rápidamente vas a incorporar todos estos nuevos conceptos. Mientras tanto, puedes ayudarte con el Glosario si algo te confunde. Challenge: prepara el entorno para trabajar con tus notebooks Parte 1. Instala Python. Asegúrate de estar conectado/a a una red segura y estable, y sigue los pasos a continuación:

Parte 2. Instala Jupyter. Los siguientes pasos son comunes a todos los sistemas operativos: Parte 3. Comprende Google Colab:

**Instala Python** 

**Google Colab** 

**Profundiza** 

Te invitamos a conocer más sobre el tema de esta bitácora.

