

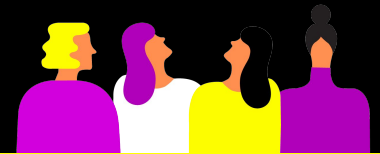
# Carrera Data Science



*Labs*  
*Proyecto Grupal*  
**#soyhenry**



## Proyecto Grupal



## Taxis NYC & Weather

Datos históricos de viajes en taxis de la ciudad de Nueva York y una API del clima

Diversos KPIs y una serie de correlaciones entre viajes y clima



### KPIs

- Días de la semana con más viajes
- Barrios con mayor participación
- Correlación entre frío/calor y viajes
- Analytics sobre viajes/pasajeros/montos

## Olist

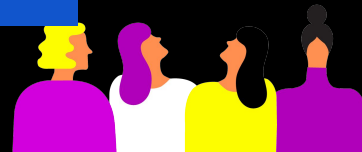
Datos históricos de compras y envíos de una de las empresas más grandes de E-commerce de Brasil

Diversos KPIs y distintas correlaciones entre compras y estadíos pandémicos



### KPIs

- Performance del delivery
- Feedback de los productos y clientes
- Plcos de ventas
- Meses con mejor revenue



## NYC tránsito y siniestralidad vial

Datos históricos de siniestralidad vial en la Ciudad de Nueva York. Se dispone de una tabla principal y varias secundarias para incorporar más información.

Además se pueden investigar API's, por ejemplo de clima.

Diversos KPIs y correlaciones entre siniestralidad y condiciones climáticas, horas pico del día, días particulares de la semana, etc.



Keywords de ejemplo:

- Días de la semana con mayor siniestralidad
- Tipo de transporte que más cantidad de accidentes genera
- Puntos de mayor concentración de accidentes

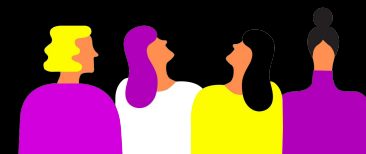
## Consumo energético y generación CO2

Datos recopilados en 3 datasets en formato .csv de diversas fuentes de generación de energía eléctrica por fuente, por planta generadora, por generación de CO2, datos de población, PBI por país y por año. Proponemos investigar API's y datasets adicionales para complementar el análisis.



Keywords de ejemplo

- Intensidad de carbono
- Huella de carbono
- Energías renovables



# LABS: Proyecto Grupal

## Objetivo Final



### Data Ingest

Dado una cantidad de datasets y API, poder obtener la estructura y datos

- 🔧 Docker
- 🔧 Python (pandas, numpy)
- 🔧 MinIO Local, S3 compatible object-storage

### Data Lake Storage

Almacenar los datos con un mínimo de limpieza y normalización

- 🔧 Docker
- 🔧 Python (pandas, numpy)
- 🔧 Nifi

### Data process

Mediante distintas técnicas y algoritmos vamos a proceder a actualizar nuestro sistema de almacenamiento de dato estructurado

- 🔧 Docker
- 🔧 Python (pandas, numpy)
- 🔧 Airflow
- 🔧 SQL

### Data Warehouse

Sistema de almacenamiento de datos estructurados, sobre el cual la organización va a obtener sus datos para la toma de decisiones

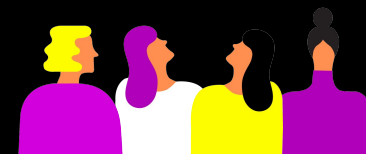
- 🔧 Docker
- 🔧 Python (pandas, numpy)
- 🔧 Airflow
- 🔧 SQL



### Data Analytics

Mediante reportes y visualizaciones vamos a facilitar la toma de decisiones

- 🔧 Python (pandas, numpy)
- 🔧 Airflow
- 🔧 SQL
- 🔧 PowerBI



# LABS: Proyecto Grupal

## Cronograma



	W1 - Data Ingest			W2 - Data Process			W3 - Data Analytics			W4 - Demo Final		
	Daily	Weekly	Demo	Daily	Weekly	Demo	Daily	Weekly	Demo	Daily	Weekly	Demo
Lunes	✓	✗	✗	✓	✗	✗	✓	✗	✗	✓	✗	✗
Martes	✓	✗	✗	✓	✗	✗	✓	✗	✗	✓	✗	✗
Miercoles	✓	✗	✗	✓	✗	✗	✓	✗	✗	✗	✗	✓
Jueves	✓	✗	✗	✓	✗	✗	✓	✗	✗	✗	✗	✓
Viernes	✗	✓	✓	✗	✓	✓	✗	✓	✓	✗	✗	✓
Objetivo	Entender el alcance del proyecto y los datasets propuestos. Diseñar una solución y entregables.			Ingesta total de datos en un Data Lake local Diseño y creación del DW Creación de los Pipelines que alimentan el DW			Diseño y creación de reportes y visualizaciones KPIs a destacar Distintos niveles de presentación para distintas audiencias					

# LABS: Proyecto Grupal

## Hitos y baseline



Semana #1	Semana #2	Semana #3	Semana #4
<p><b>Puesta en marcha el proyecto</b></p> <ul style="list-style-type: none"><li>• Kickoff del proyecto</li><li>• Entendimiento de las necesidades</li><li>• Documentar alcance, objetivo y entregables</li></ul>	<p><b>Trabajando los datos</b></p> <ul style="list-style-type: none"><li>• Creación del DW</li><li>• Reglas de negocio aplicadas</li><li>• Automatizar el DW</li></ul>	<p><b>Etapas de Analytics</b></p> <ul style="list-style-type: none"><li>• Reportes</li><li>• Storytelling</li><li>• Ajustes necesarios al modelo</li></ul>	<p><b>Retoques finales y presentación</b></p> <ul style="list-style-type: none"><li>• Preparar demo por equipo</li><li>• Entregable final</li><li>• Documentación</li></ul>

# LABS: Proyecto Grupal

## Hitos - Semana #1



### Semana #1

#### Puesta en marcha el proyecto

- Kickoff del proyecto
- Entendimiento de las necesidades
- Documentar alcance, objetivo y entregables

1. Entendimiento de la situación actual

2. Objetivos

3. Alcance

4. Fuera de alcance

5. Solución propuesta - Incluir Stack tecnológico

6. Metodología de trabajo

7. Diseño detallado – Entregables

8. Equipo de trabajo – Roles y responsabilidades

9. Cronograma general



# LABS: Proyecto Grupal

## Hitos - Semana #2



---

### Semana #2

#### Trabajando los datos

- Creación del DW
- Reglas de negocio aplicadas
- Automatizar el DW

1. Diseño adecuado del Modelo

2. Documentación

3. Pipelines para alimentar el DW

4. Automatización

5. Validación de datos

# LABS: Proyecto Grupal

## Hitos - Semana #3



---

### Semana #3

#### Etapa de Analytics

- Creación del DW
- Reglas de negocio aplicadas
- Automatizar el DW

1. Diseño de Reportes/Dashboards

2. Documentación

3. Pipelines para alimentar el DW

4. Automatización

5. Validación de datos

# LABS: Proyecto Grupal

## Hitos - Semana #4



---

### Semana #4

#### Retoques finales y presentación

- Preparar demo por equipo
- Entregable final
- Documentacion

1.Prepara la demo, visualización efectiva

2. Documentación

3. Probar todo el proceso antes!!!

# LABS: Proyecto Grupal

## Baseline esperado



Semana #1	Semana #2	Semana #3	Semana #4
<p><b>Puesta en marcha del proyecto y definiciones iniciales:</b></p> <ul style="list-style-type: none"><li>• Al menos 4 KPIs</li><li>• Tecnologías a usar</li><li>• Documento de alcance del proyecto</li></ul>	<p><b>Trabajando los datos</b></p> <ul style="list-style-type: none"><li>• Datawarehouse automatizado con carga inicial. Al menos 2 tablas de hechos y 5 dimensionales</li></ul>	<p><b>Etapas de Analytics</b></p> <ul style="list-style-type: none"><li>• Carga incremental</li><li>• Dashboard y reportes</li></ul>	<p><b>Retoques finales y presentación</b></p> <ul style="list-style-type: none"><li>• La presentación debe estar dirigida a la dirección de la Compañía</li><li>• Storytelling</li></ul>
PLUS	PLUS	PLUS	PLUS
<ul style="list-style-type: none"><li>• Incrementar número de KPIs</li><li>• Planificación y estimación de esfuerzos. Diagrama Gantt.</li></ul>	<ul style="list-style-type: none"><li>• Uso de herramientas Big Data como HDFS, Hive, Spark y/o motores No-SQL</li></ul>	<ul style="list-style-type: none"><li>• Implementar modelo de Machine Learning</li></ul>	<ul style="list-style-type: none"><li>• Implementar un reporte con visualización geográfica</li></ul>

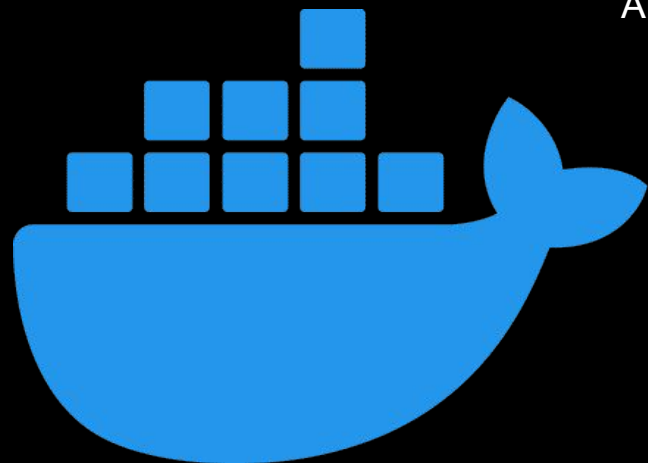
# LABS: Proyecto Grupal

## Docker para trabajar el PF



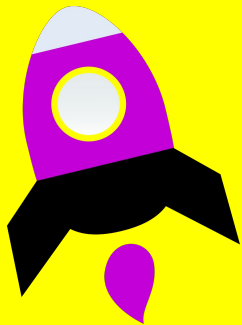
<https://github.com/sercasti/datalaketools>

Alternativa <https://github.com/Marcel-Jan/docker-hadoop-spark>



docker®

# Q&A



**#soyhenry**



# Muchas Gracias

**#soyhenry**

