



NOMBRE DEL PROYECTO	MENTOR	PATROCINADOR
SMART CHOICE ANALYTICS	EMILIO SANTANDER	HENRY

	EXPECTED START DATE	EXPECTED COMPLETION
Proyecto Final Data Science Henry PT-07	29/02/2024	07/06/2024

Equipo del Proyecto

Rol	Name	e-mail	Name-Abr	Git Hub
Machine Learning	Diego Osorio	dosoriofc@gmail.com	DO	dosoriofc
Data Engineer	David Marimón	david.neko26@gmail.com	DM	DaAnMaGi
Data Engineer	Salomon Orozco	Salomonorozcojaramillo@gmail.com	SO	SaloLL
Data Analytics	Marcela Correal	mcorreal@gmail.com	MC	MarceCorreal
Data Analytics	Juan Garate	garatejb@gmail.com	JG	Batxa

Links asociados

- Documentos internos
 - Repositorio del proyecto https://github.com/Batxa/DS_ProjectFinal.git
 - Follow-up de tareas: <https://trello.com/b/sH9ofad9/tpfinal>
 - Cronograma: <https://app.teamgantt.com/projects/gantt?ids=3939144>
- Documentos externos (Soy Henry)
 - Presentación y estructura del proyecto: https://docs.google.com/presentation/d/17WLH3D5cgEbNA3XIKYKMEi20Q9pbo6kOXmD74nRhQal/edit#slide=id.g1f3119b515c_0_841
 - Enunciado del proyecto (README): https://github.com/soyHenry/PF_DS/blob/FULL-TIME/README.md
 - Diccionario de datasets originales" https://docs.google.com/document/d/1ASLMGAgviiCAtaP1UJlflpmBCXuSTHQQGWdQMn6_2l/edit
 - Detalles de estructura, hitos y entregables: https://docs.google.com/document/d/e/2PACX-1vRtTsN_N3Z0DTLbh_-Xw2OxhOWeV5jmTISRUNzTBpWM9mTnxST03674UheR4f0hfULc2v4_sW3lgDTv/pub

- Rúbrica de evaluación:
https://docs.google.com/document/d/1tBuh1LSCmvQB5Wd7-Cj4jj_o5zLk8vtBQFtDhF8oeSY/edit#heading=h.okiecvgw2s1c

Metodología de trabajo

Se utilizará la metodología Scrum-Agile, en la cual se aplicarán los siguientes conceptos:

- Iteraciones cortas (Sprints):
 - El trabajo se divide en 3 sprints de 2 semanas de duración cada uno.
 - Al final de cada sprint se debe obtener uno o más entregables.
- Reuniones diarias (Daily Stand Ups):
 - Reuniones diarias de 45 minutos donde el equipo comparte actualizaciones sobre lo que han hecho desde la última reunión, lo que planean hacer hoy y cualquier obstáculo que estén enfrentando.
 - El resultado de cada reunión será la priorización de tareas y la indicación de los próximos pasos de cada integrante
- Colaboración con el cliente:
 - Implica trabajar en estrecha colaboración con el cliente (mentor) bajo una comunicación abierta y continua, en función de comprender sus necesidades, expectativas y cualquier cambio requerido.
 - El mayor beneficio es que se obtiene un feedback temprano y regular, garantizando que el producto final satisfaga las expectativas.
- Demostraciones y retrospectivas:
 - Al final de cada sprint, se lleva a cabo una demostración del trabajo completado y una retrospectiva para identificar qué se hizo bien, qué se podría mejorar y cómo.
- Entrega continua:
 - Se busca la entrega de software en incrementos pequeños y frecuentes, en lugar de esperar hasta el final del proyecto para entregarlo todo de una vez.
- Equipos multifuncionales y autoorganizados:
 - Los equipos Agile son multifuncionales, lo que significa que incluyen todas las habilidades necesarias para completar el trabajo, y son autoorganizados, lo que les permite tomar decisiones y gestionar su propio trabajo.

Generalidades del Proyecto

Cliente	El proyecto está dirigido a propietarios y directivos de restaurantes y bares de Estados Unidos, interesados en obtener información clave sobre las tendencias e intereses de los clientes del sector, con el propósito de que ésta sirva como input para sus proyectos y estrategias que mejoren el desempeño de su marca.
Alcance	Los datos que involucra este proyecto se limitan a los restaurantes y bares ubicados en el territorio de los Estados Unidos, de todas las categorías desde el año 2010 al año 2021

	<p>Los mismos se disponibilizarán en un Data WareHouse que podrá ser consultada en cualquier momento por el cliente.</p> <p>Los reviews son proporcionados por la plataforma Yelp, y los encontrados en el sistema Google Maps.</p>
Objetivo 1	<p>Desarrollar una herramienta de gestión de la información (Dashboard) que ayude a los directivos de restaurantes a mejorar la elaboración de estrategias o campañas de marketing y los servicios ofrecidos a sus clientes a través de las siguientes funcionalidades:</p> <ul style="list-style-type: none"> a. Segmentación de clientes: Hacer uso de técnicas de clustering para segmentar a los clientes en grupos homogéneos, en función de sus patrones de consumo y comportamiento, contribuyendo a personalizar ofertas y estrategias de marketing. b. Análisis de sentimientos en reseñas: Aplicar técnicas de procesamiento de lenguaje natural (NLP) y análisis de sentimientos para comprender las opiniones de los usuarios sobre el negocio, extrayendo información valiosa sobre los intereses y gustos de los clientes. c. Detección de tendencias en la industria: Emplear técnicas de análisis de datos y series temporales para identificar tendencias emergentes en la industria de restaurantes, bares y negocios a nivel local y nacional, contribuyendo a la adaptación oportuna del negocio a los cambios en la demanda del mercado. <p>La herramienta permitirá al usuario ejecutar filtros de acuerdo a establecimientos, lugares, categorías, reviews y calificaciones, lo que permitirá sacar conclusiones de valor para el desarrollo de estrategias.</p>
Objetivo 2	<p>Presentar un análisis del mercado de restaurantes y bares estadounidense basados en las reseñas de clientes de las plataformas de Yelp y Google Maps, identificando e ilustrando las características más relevantes del mercado, mediante el análisis de las relaciones entre las variables que generan un mayor impacto en las evaluaciones de los establecimientos.</p> <p>El estudio, basado en la utilización de la herramienta, incluirá elementos tales como:</p> <ul style="list-style-type: none"> a. Características de los establecimientos tales como: Categoría, Atributos/Misceláneos, Nivel de precios, reviews, Rating/Stars asignados b. Cantidad de comentarios positivos/neutros/negativos (obtenido del análisis de sentimiento) c. Características de los usuarios tales como: número de reseñas suministradas, número de años registrado en Yelp, número de amigos en Yelp d. Ranking de Restaurantes y bares por Comentarios positivos (obtenido del análisis de sentimiento, con número de reseñas como segundo criterio de ordenamiento) e. Ranking de Restaurantes por Comentarios Negativos (obtenido del análisis de sentimiento, con número de reseñas como segundo criterio de ordenamiento)

Entregables Finales	Herramienta de Análisis de Datos: Dashboard que brinda la posibilidad de filtrar por cada uno de los segmentos definidos, y así mismo, visualizar dinámicamente el movimiento que los indicadores presentan, de acuerdo con los intereses del estudio particular.
	Análisis del sector gastronómico en base a los datos que se analizan con la herramienta a entregar.
	<p>Sistema de Clasificación desarrollado con herramientas de <i>Machine Learning</i>, que permita el reconocimiento y segmentación de grupos de clientes.</p> <p>Los cuales serán entregados de la siguiente manera:</p> <p>Sprint 1:</p> <ol style="list-style-type: none"> 1. Charter del proyecto (incluye fundamentos del proyecto, stack tecnológico y flujo de trabajo, Alcance, objetivo, entregables). 2. Exploración de los datos (EDA). <p>Sprint 2:</p> <ol style="list-style-type: none"> 3. Reporte de Data Analytics y producto Machine Learning. <p>Sprint 3:</p> <ol style="list-style-type: none"> 4. Dashboard final. 5. Producto final de Machine Learning.
Sprint 1 EDA (Exploratory Data Analysis)	<p>Para la realización del EDA, se realiza en forma previa la fase de ETL (Extraction, transformation and loading of data) en forma separada, debido a la gran cantidad de información. La estructura y contenido del EDA es el siguiente:</p> <ul style="list-style-type: none"> * Importación de librerías * Carga de datos * Pre-procesamiento de datos: gestión de tipos de datos, valores duplicados y nulos, gestión de características * Análisis de datos: incluye series de tiempo, distribuciones y formatos diversos para analizar la evolución a través del tiempo de características * Conclusiones y resumen de hallazgos
Sprint 1 KPIs	<p>Se incluirán en el Dashboard 3 indicadores, diseñados para apoyar el diseño de las estrategias de mejora de los restaurante y bares. Los mismos estarán graficados dinámicamente en el Dashboard, herramienta que permitirá iterar gráficamente suposiciones de comportamiento.</p> <p>Los 3 KPIs son los siguientes:</p>
KPI 1	Aumentar en un (10%) Número de Calificaciones con valores mayores a (3.5), en el último trimestre
KPI 2	Aumentar en un (10%) Número de Reseñas consideradas como Positivas en el último trimestre
KPI 3	Aumentar en un 10% las ventas del último trimestre

Esquema de Trabajo

El proyecto consta de un total de 3 bloques o sprints con las siguientes características:

- Sprint 1: Puesta en marcha
 - Duración: 2 semanas
 - Actividad principal: Contextualización de situación y propuesta de solución
- Sprint 2: Data Engineering
 - Duración: 2 semanas
 - Actividad principal: Implementación de estructura de datos
- Sprint 3: Data Analytics y Machine Learning
 - Duración: 2 semanas
 - Actividad principal: Desarrollo final de la herramienta de análisis de la información y el producto de *Machine Learning*.

Stack Tecnológico	
Lenguaje Principal de programación	Python
Gestor Base de datos	SQL
Repositorio para alojar código	Git Hub
Herramienta de seguimiento del Proyecto	Trello
Plataforma de lanzamiento	Visual Code Studio
Bibliotecas de Python	Pandas, Numpy, LTK, SCikitl Learn
Servicio de Nube	GCP - Fabric

Datasets a utilizar

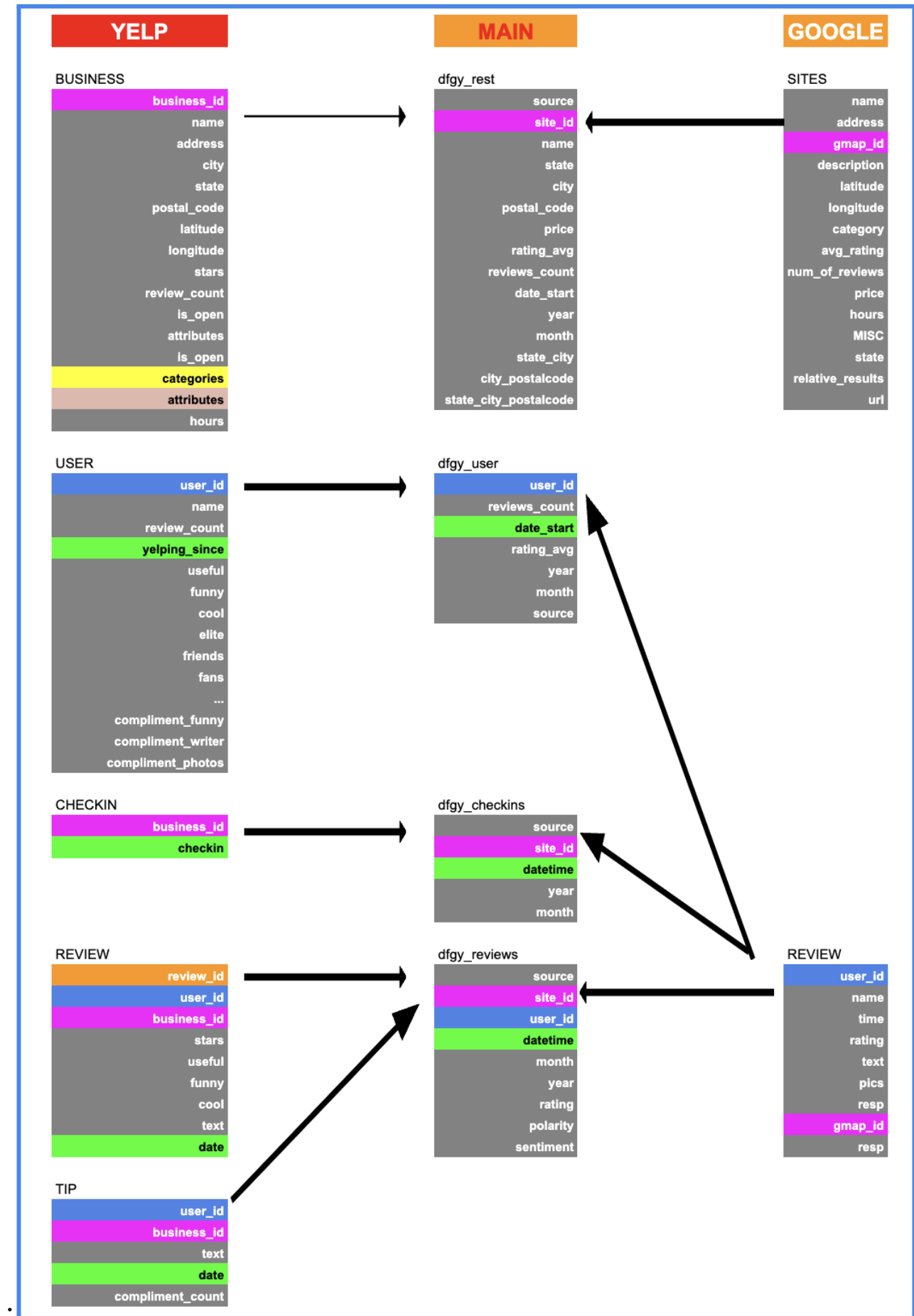
- Fuentes de datos obligatorios
 - Datasets de YELP
 - Business: contienen la información del comercio
 - Checkin: contiene registros de check-in o visitas de clientes
 - User: contiene la data propia del usuario
 - Tip: son consejos o sugerencias rápidas por parte del usuario
 - Review: contiene las reseñas completas
 - Datasets de GOOGLE
 - Sites: contiene información de los sitios
 - Reviews: contiene reseñas de usuarios
- Fuentes de datos adicionales
 - Códigos postales de USA
 - Datos demográficos

- Datos de crecimiento económico
- <https://data.gov/>

Alcance de la información a utilizar:

- Alcance de rubros: restaurantes y bares
- Alcance geográfico: Estados Unidos
- Alcance temporal: 2010 en adelante, la información llega hasta el 2021
- Sesgos posibles: la información contiene datos del año 2020 siendo el covid 19 un factor de alto impacto. A efectos del EDA se incluirá esta información, y luego se definirá si se utiliza o no para las proyecciones del producto final

Flujo y estructura de datos:



Plan y cronograma de trabajo (Esquema Gantt)

