

Universidad Panamericana  
Maestría en Ciencia de Datos  
Datos Masivos

Proyecto Final: *Monitoreo Inteligente de Tráfico Marítimo con  
Big Data AIS*

Enrique Ulises Báez Gómez Tagle, Luis Alejandro Guillén Alvarez

26 de septiembre de 2025

## Índice

<b>1</b>	<b>Resumen ejecutivo</b>	<b>2</b>
1.1	Introducción . . . . .	2
<b>2</b>	<b>Visión general del desarrollo</b>	<b>2</b>
2.1	Solución Actual / visión general . . . . .	2
2.2	Limitaciones actuales de la solución . . . . .	3
2.3	Propósito, uso y alcance de la herramienta . . . . .	3
<b>3</b>	<b>Revisión y uso de datos</b>	<b>5</b>
3.1	Orígenes y control de datos . . . . .	5
3.2	Preparación de datos . . . . .	5
3.3	Limpieza y tratamiento de datos . . . . .	6
3.4	Limitaciones de los datos . . . . .	7
3.5	Integridad de los datos . . . . .	8
<b>4</b>	<b>Proceso de desarrollo</b>	<b>9</b>
4.1	Metodología . . . . .	9
4.2	Pruebas . . . . .	9
<b>5</b>	<b>Resultados y conclusiones</b>	<b>11</b>
5.1	Conclusión general . . . . .	14
<b>6</b>	<b>Código utilizado</b>	<b>15</b>
6.1	Link al repositorio con código fuente y salidas correspondientes . . . . .	15

# 1. Resumen ejecutivo

## 1.1. Introducción

El tráfico marítimo global genera millones de registros de posicionamiento a través del Sistema de Identificación Automática (AIS), que transmite automáticamente la identidad, posición, velocidad y otros datos de los buques. Originalmente implementado para mejorar la seguridad de la navegación, hoy es crucial para la gestión del tráfico marítimo, la conciencia del entorno y operaciones de búsqueda y rescate.

Este proyecto propone aprovechar Big Data y aprendizaje automático (ML) para analizar una base de datos masiva de mensajes AIS y ofrecer una solución innovadora a la pregunta: **¿Cómo detectar y entender comportamientos anómalos en el tráfico marítimo para mejorar la seguridad y la eficiencia?**

Responder esta pregunta implica analizar patrones normales de navegación y descubrir desviaciones significativas. En esencia:

- **¿Qué estamos haciendo?** Diseñamos un sistema de análisis que procesa grandes volúmenes de datos AIS para identificar anomalías en el comportamiento de los buques (posiciones fuera de lugar, velocidades inusuales, maniobras erráticas, etc.) y extraer patrones útiles sobre la operación de diferentes tipos de embarcaciones.
- **¿Para qué lo hacemos?** Para mejorar la seguridad marítima, la gestión del tráfico y la toma de decisiones, ofreciendo alertas tempranas de potenciales riesgos (colisiones, actividades ilícitas o fallos) y conocimiento profundo a autoridades y empresas navieras.
- **¿Cómo lo hacemos?** Empleando herramientas de Big Data de la Suite de Google Cloud (Storage + DataProc + BigQuery) para procesar datos geospaciales masivos en poco tiempo, complementado con modelos de ML que aprenden patrones habituales y detectan comportamientos que se apartan de lo normal.
- **¿A quién beneficia?** A organismos de seguridad marítima (marinas, guardacostas), a empresas navieras optimizando rutas y monitoreo de flotas, a aseguradoras evaluando riesgos operativos e incluso a investigadores del medio marino en estudios ecológicos.

Finalmente, los resultados se integraron en un *dashboard* interactivo con visualizaciones (mapas geográficos, gráficos de tendencias y rankings) que facilitan la interpretación y la toma de decisiones.

## 2. Visión general del desarrollo

### 2.1. Solución Actual / visión general

La solución actual se implementa de extremo a extremo sobre la nube de Google Cloud. El flujo completo es el siguiente:

- Una máquina virtual en Google Compute Engine se encarga de **scrapear los datos AIS** y cargarlos en Google Cloud Storage.
- En Cloud Storage se organizan dos capas diferenciadas: *raw* y *curated*, con prefijos de carpeta y particiones por mes (YYYY-MM=).
- Un clúster de Google Cloud DataProc (perfil *high-memory*; véase Sección 4.1) ejecuta dos jobs, habitando procesamientos por series de meses con menos derrames a disco y menor tiempo total:
  - Job **raw**: descomprime los archivos obtenidos y organiza los datos en particiones por año y mes, y finalmente los guarda en formato Parquet en Google Cloud Storage.
  - Job **curated**: lee la capa cruda, aplica transformaciones iniciales y genera la capa refinada lista para análisis.

- Una Cloud Function crea el **dataset y la tabla en BigQuery**, cargando la información desde la capa curated.
- En BigQuery se centraliza la explotación de datos mediante consultas SQL optimizadas, generando los datasets finales que alimentan las visualizaciones.
- El **dashboard de Streamlit**, desplegado en la misma VM, consume los resultados de BigQuery y presenta mapas, gráficos de tendencias y rankings de manera interactiva.

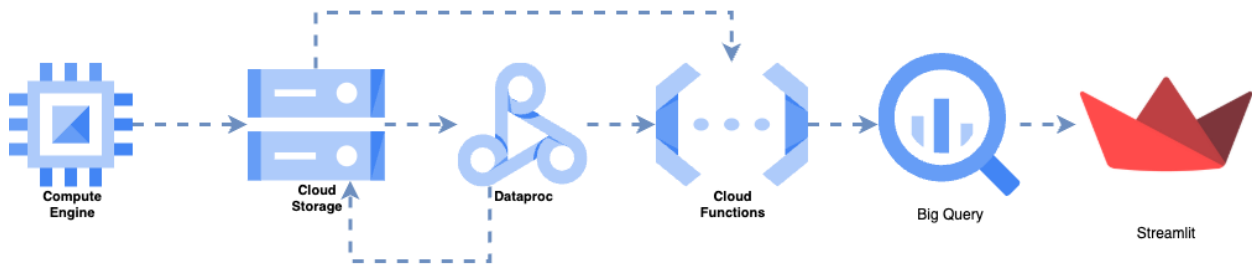


Figura 1: Arquitectura actual de la solución: flujo end-to-end en Google Cloud.

## 2.2. Limitaciones actuales de la solución

- **Sin orquestación event-driven.** Los disparadores son manuales o programados; la activación automática por eventos de GCS/PubSub o Workflows aún no está implementada. Migrar a un esquema event-driven reduciría intervención manual y tiempos de ciclo.
- **No tiempo real.** El pipeline es *batch*; no hay ingesta en streaming ni SLA de latencia. El dashboard refleja cortes al cierre de cada corrida.
- **Evolución de esquema limitada.** La Cloud Function crea tablas y puede autodetectar esquema, pero no hay un proceso formal de migraciones/versionado ni *backfills* automatizados.
- **Pruebas y monitoreo.** Podrían generarse alertas proactivas (Cloud Monitoring/Error Reporting) ante fallas o derivas.

## 2.3. Propósito, uso y alcance de la herramienta

**Propósito general.** La herramienta tiene como propósito transformar los mensajes AIS masivos en *inteligencia operable* para seguridad y gestión del tráfico marítimo. Provee un flujo reproducible de ingesta, depuración, análisis y detección de anomalías que permita identificar comportamientos de riesgo (velocidades inusuales, cambios de rumbo bruscos, calados atípicos, derivas de patrones por zona/tipo de buque) y entregar hallazgos mediante un *dashboard* interactivo.

### Objetivos específicos.

- **Monitoreo descriptivo y exploratorio:** caracterizar el comportamiento típico por tipo de embarcación (*VesselTypeName*), zona (*geohash9*) y periodo (hora/día), incluyendo estados de navegación y métricas de velocidad/rumbo.
- **Detección temprana de anomalías:** señalar observaciones que se desvían significativamente de los patrones normales (p.ej., *z-scores* de calado, picos de velocidad, virajes abruptos) y, a nivel temporal, detectar valores atípicos con modelos ARIMA\_PLUS en BigQuery ML.
- **Soporte a decisiones:** priorizar buques, zonas y ventanas temporales que requieren atención, entregando listas y visualizaciones accionables para autoridades y operadores.

- **Escalabilidad y reproducibilidad:** estandarizar el pipeline en Google Cloud (Storage, DataProc, BigQuery, Cloud Functions) para procesar volúmenes elevados con costos y tiempos controlados.

#### Usuarios objetivo y decisiones habilitadas.

- **Autoridades marítimas/guardacostas:** focalización de patrullajes y verificación de eventos inusuales por zona-tiempo-buque.
- **Navieras y operadores de flota:** supervisión de cumplimiento operativo, eficiencia de rutas y detección de riesgos operativos.
- **Aseguradoras y analistas de riesgo:** evidencia cuantitativa para evaluación de siniestralidad y *pricing* técnico.
- **Investigación/medio marino:** análisis de patrones de movilidad y su interacción con áreas sensibles.

**Casos de uso cubiertos en el *dashboard*.** Cada módulo se respalda con consultas SQL parametrizables en BigQuery:

- **Calado anómalo (*per buque y tipo*):** cálculo de *z-score* sobre el calado medio por MMSI dentro de su categoría, para resaltar posibles inconsistencias operativas o de reporte.
- **Cambios de dirección bruscos (*virajes*):** detección de deltas de rumbo (COG) sobre series temporales por MMSI, agregados por *geohash9*, con umbrales configurables de severidad.
- **Velocidades inusuales:** comparación del *máximo* de SOG contra percentiles de referencia por buque/tipo para evidenciar excesos.
- **Variabilidad de velocidad y rumbo:** STDDEV y AVG de SOG/COG por tipo de embarcación como indicador de estabilidad operacional.
- **Correlaciones estructurales:** relaciones Length{Width (eslora-manga) por clase de buque y correlaciones entre métricas operativas (p.ej., SOG vs. COG) con filtros de tamaño muestral.
- **Resumen de estados de navegación:** conteos, embarcaciones únicas y promedios (SOG, calado) por NavStatusName.
- **Patrones semanales:** velocidad promedio y estado de navegación más frecuente por día de la semana y tipo de buque.
- **Ubicación y trazas recientes (*geohash9*):** listados y mapas de últimas posiciones por filtros de fecha, tipo de buque o MMSI.
- **Detección de anomalías temporales con ML:** entrenamiento y consulta de modelos ARIMA\_PLUS por serie (*geohash9*, VesselTypeName o MMSI) para marcar valores atípicos con umbrales de probabilidad ajustables.

#### Entradas, procesamiento y salidas.

- **Entradas:** archivos AIS 2024 (NOAA) descargados y versionados en GCS (*raw/curated*), con partición temporal.
- **Procesamiento:** DataProc (Spark) para descompresión, partición y normalización; BigQuery para *modelado analítico*, agregaciones y ML.
- **Salidas:** tablas refinadas y *views* listas para consumo, modelos ML en *dataset* dedicado y un *dashboard* Streamlit con mapas, tendencias y rankings exportables.

### Alcance (in/out-of-scope).

- **Incluye:** análisis *batch*, filtros por ventana temporal y región, métricas descriptivas, detección de anomalías no supervisada (estadística y ARIMA\_PLUS), visualizaciones interactivas y *ranking* de eventos/buques.
- **No incluye (versión actual):** ingesta *streaming* tiempo real, orquestación *event-driven*, *backfills* automáticos de esquemas, clasificación supervisada de ilícitos, reconstrucción de trayectorias de alta fidelidad ni generación de alertas operativas en caliente.

### Supuestos y restricciones.

- Calidad y cobertura del AIS dependen de la fuente; pueden existir vacíos, duplicados o reportes erróneos.
- El desempeño está condicionado por la configuración del clúster DataProc y los límites/cuotas de BigQuery; se privilegia costo-eficiencia sobre latencia.
- La zonificación espacial usa `geohash9` como índice geoespacial operacional.

## 3. Revisión y uso de datos

### 3.1. Orígenes y control de datos

Los datos utilizados provienen del portal oficial de la NOAA (<https://coast.noaa.gov/htdata/CMSP/AISDataHandler/2024/index.html>), el cual publica archivos comprimidos (.zip) con los mensajes AIS recolectados a lo largo del año 2024. Para garantizar la trazabilidad y el control de los datos, se desarrolló un script en Python (`scraper.py`) que automatiza la descarga de todos los archivos disponibles, verificando duplicados mediante el tamaño de los ficheros y evitando descargas innecesarias. De esta manera se asegura que la totalidad de los datos originales publicados por la NOAA sean capturados sin alteración.

### 3.2. Preparación de datos

En esta sección se describe el **qué** produce la preparación de datos (entradas, salidas y *esquema*); la orquestación y operación del pipeline (el **cómo**) se documenta en la Sección 4.1.

La preparación de los datos comienza en una máquina virtual de Google Compute Engine, donde se ejecuta el script de *scraping*. Este proceso descarga en paralelo (con *threads*) los archivos .zip hacia un directorio local (`ais_2024`) y posteriormente los traslada a Google Cloud Storage en la capa **raw**.

Dado el objetivo de tener un conjunto uniforme, trazable y listo para análisis, con: (i) esquema estandarizado, (ii) claves y derivadas temporales, (iii) coordenadas y velocidades normalizadas, y (iv) datos particionados por *ym* (AAAA-MM).

#### Entradas → salidas por etapa y la salida de cada una.

1. **RAW INGESTION:** entrada = lotes .zip de NOAA; salida = archivos Parquet en **raw** *particionados por ym*, con trazabilidad mínima (archivo fuente y timestamp de ingesta).
2. **CURATED TRANSFORMATIONS:** entrada = **raw**; salida = capa **curated** con tipos normalizados, derivados temporales y deduplicación por `MMSI+BaseDateTime`. Las reglas de limpieza y validaciones se detallan en la Sección 3.3.

**Modelo de datos de la capa *curated*.** El Cuadro 1 presenta los campos principales que expone la capa analítica.

Cuadro 1: Esquema de la capa *curated*.

Campo	Tipo	Descripción
NavStatusInt	INTEGER	Representación entera del estado de navegación del buque.
MMSI	STRING	Número de Identidad del Servicio Móvil Marítimo del buque.
BaseDateTime	TIMESTAMP	Marca temporal que indica la fecha y hora del mensaje AIS.
LAT	FLOAT	Latitud de la posición del buque.
LON	FLOAT	Longitud de la posición del buque.
SOG	FLOAT	Velocidad sobre el fondo del buque en nudos.
COG	FLOAT	Rumbo sobre el fondo del buque en grados.
Heading	FLOAT	Rumbo de la proa del buque en grados.
VesselName	STRING	Nombre del buque.
IMO	STRING	Número de la Organización Marítima Internacional del buque.
CallSign	STRING	Indicativo de llamada por radio del buque.
VesselType	STRING	Tipo de buque.
Status	STRING	Estado del buque.
Length	FLOAT	Eslora del buque.
Width	FLOAT	Manga del buque.
Draft	FLOAT	Calado del buque.
Cargo	STRING	Tipo de carga que transporta el buque.
TransceiverClass	STRING	Clase del transpondedor AIS.
_source_file	STRING	Nombre del archivo fuente desde el cual se ingirieron los datos.
_ingest_ts	TIMESTAMP	Marca temporal que indica cuándo se ingirieron los datos en la tabla.
ymd	STRING	Fecha en formato año-mes-día.
VesselTypeInt	INTEGER	Representación entera del tipo de buque.
VesselTypeCode	STRING	Código que representa el tipo de buque.
VesselTypeName	STRING	Nombre del tipo de buque.
VesselTypeClass	STRING	Clasificación del tipo de buque.
NavStatusName	STRING	Nombre del estado de navegación del buque.
date	DATE	Fecha en que se registró el mensaje AIS.
hour	INTEGER	Hora del día en que se registró el mensaje AIS.
dow	STRING	Día de la semana en que se registró el mensaje AIS.
week	INTEGER	Semana del año en que se registró el mensaje AIS.
month	INTEGER	Mes en que se registró el mensaje AIS.
quarter	INTEGER	Trimestre del año en que se registró el mensaje AIS.
SOG_ms	FLOAT	Velocidad sobre el fondo del buque en metros por segundo.
geohash9	STRING	Representación geohash de la posición del buque con precisión de 9 caracteres.

Una vez procesado, el resultado de este pipeline se publica en BigQuery como tabla particionada por *ym*, optimizada para consulta y consumo por el *dashboard*. Los detalles de automatización de creación/carga (Cloud Function, idempotencia y autodetección de esquema) se abordan en la Sección 4.1.

### 3.3. Limpieza y tratamiento de datos

En esta sección se describen de manera detallada las transformaciones y criterios aplicados para convertir los datos AIS en una capa **curated** funcional para el análisis. Se listan los pasos con su *porqué* operativo o estadístico.

#### Normalizaciones y tipos.

- **Casteo y estandarización:** se convierten columnas numéricas a tipos adecuados (*double/int*) y se normalizan textos (*trim*, mayúsculas controladas) para evitar variabilidad en joins/agrupaciones.
- **Timestamps en UTC:** *BaseDateTime* se parsea a tipo *timestamp* con zona UTC para análisis temporales coherentes; sobre esta base se derivan *date*, *hour*, *week*, *month*, *quarter* y una columna *ym* (AAAA-MM) para particionado.

#### Coordenadas y límites físicos.

- **LAT/LON válidos:** se recortan latitudes al rango  $[-90, 90]$  y longitudes se *envuelven* a  $[-180, 180]$  (corrigiendo valores fuera de dominio por errores/reporte). Se descartan registros con coordenadas nulas tras la corrección.

- **Redondeo razonable:** se redondean LAT/LON a 5 decimales para homogenizar precisión sin perder resolución náutica relevante.

#### Velocidad, rumbos y estados.

- **SOG (Speed Over Ground):** se acota a [0, 70] nudos para eliminar outliers físicos/telemetría corrupta; adicionalmente se calcula **SOG<sub>ms</sub>** en m/s para modelos.
- **COG/Heading:** se normaliza módulo 360°. El valor **Heading=511** (no disponible) se marca como nulo para evitar sesgos en métricas de rumbo.
- **NavStatus:** se genera **NavStatusInt** a partir de **Status** y se une contra un catálogo para **NavStatusName**; códigos fuera de 0–15 se etiquetan como *Unknown* y faltantes como *Not reported*.

#### Dimensiones y calado.

- **Length/Width:** se acotan a rangos plausibles (1–450 m y 1–70 m, respectivamente) para descartar errores de orden de magnitud.
- **Draft:** se limita a [0, 25] m; valores negativos o excesivos suelen provenir de carga mal reportada o errores de parsing.

#### Identificadores y texto.

- **MMSI:** se extrae la porción numérica válida (1–9 dígitos) y se rellena a 9 con ceros a la izquierda para formato consistente.
- **IMO/CallSign/VesselName:** se normalizan espacios y caracteres no alfanuméricos; se limpian IMOs genéricos (p.ej., IM00000000, 0, vacío) para evitar falsas uniones.

#### Clasificación de tipo de buque.

- **VesselType:** se derivan **VesselTypeInt**, **VesselTypeName** (mediante mapa IMO) y **VesselTypeClass** (WIG, Cargo, Tanker, Passenger, etc.). Esto permite comparar métricas por clase homogénea y soportar reglas dependientes del tipo.

#### Enriquecimiento geoespacial.

- **Geohash9:** se calcula un **geohash9** distribuido (UDF con **pygeohash**) para agregación espacial eficiente y enlaces a capas geográficas (vallas, puertos, TSS) en análisis posteriores.

#### Calidad e idempotencia.

- **Desduplicación:** se eliminan duplicados por clave compuesta (**MMSI**, **BaseDateTime**) para asegurar series temporales limpias por embarcación.
- **Validación de conteos:** en *raw*, se comparan conteos diarios (CSV vs Parquet) y se aborta si no coinciden; en *curated*, se escriben marcadores de éxito por partición (**\_markers/ym=.../\_SUCCESS**) para reanudación y ejecución incremental sin sobrescribir trabajo correcto.

### 3.4. Limitaciones de los datos

#### Limitaciones de origen y cobertura.

- **Cobertura heterogénea:** la recepción varía por densidad de estaciones y condiciones atmosféricas; en zonas costeras y puertos existe mejor cobertura que en mar abierto.

- **Congestión y pérdida de mensajes:** en áreas de alto tráfico se producen colisiones de ranuras/transmisiones y *dropouts*; los conteos absolutos pueden estar subestimados.
- **Latencia y vacíos temporales:** hay periodos con reportes ausentes o irregulares por buque; no deben interpretarse como inmovilidad o “apagado” sin evidencia adicional.

#### Limitaciones espaciotemporales y cinemáticas.

- **Posición y tiempo:** `BaseDateTime` puede no estar perfectamente sincronizado; posiciones pueden presentar *jitter* o saltos por errores GNSS.
- **Velocidad y rumbos:** `SOG/COG` contienen ruido y picos espurios; `Heading=511` indica “no disponible” y se trata como nulo.

#### Sesgos analíticos y de modelado.

#### Implicaciones para la interpretación.

- **Anomalía  $\neq$  ilícito:** una alerta es un *candidato a revisión*, no evidencia de conducta indebida por sí misma.
- **Comparaciones geográficas:** los conteos no son estrictamente comparables entre regiones con distinta cobertura; preferir tasas relativas o métricas normalizadas.
- **Temporalidad:** el análisis es *batch*; el *dashboard* refleja el último corte procesado, no tiempo real.

### 3.5. Integridad de los datos

La integridad se asegura con **controles preventivos y detectivos** desde la ingesta hasta el consumo. Objetivos: (i) preservar *fidelidad* a la fuente (NOAA), (ii) garantizar *validez/consistencia* de esquema y valores, (iii) habilitar *trazabilidad e idempotencia*, y (iv) exponer *métricas de calidad* operables.

#### Principios y dimensiones.

- **Compleitud:** cobertura esperada por periodo (*ym*).
- **Unicidad:** clave (`MMSI`, `BaseDateTime`) sin duplicados.
- **Validez:** rangos físicos plausibles (`LAT/LON`, `SOG`, `COG`, `Length/Width`, `Draft`).
- **Consistencia:** Catálogos coherentes.
- **Actualidad:** disponibilidad oportuna en `BigQuery`.
- **Trazabilidad:** `_source_file` y `_ingest_ts`.

#### Controles por etapa.

- **Ingesta (`GCE`  $\rightarrow$  `GCS/raw`):** detección de duplicados por nombre/tamaño; estructura `raw/ym=AAAA-MM/`.
- **Estandarización (`DataProc`):** normalización (`COG/Heading` mod 360, `SOG` cap, `LAT/LON` recorte), dedup (`MMSI`, `BaseDateTime`), derivadas y `geohash9`.
- **Carga analítica (`CF` + `BigQuery`):** partición por `ym`, marcadores de *SUCCESS* e idempotencia.

#### Linaje, reproducibilidad e idempotencia.

- **Linaje:** cada fila conserva archivo de origen y timestamp de ingesta.
- **Reproducibilidad:** transformaciones declarativas; particiones `ym` inmutables tras *SUCCESS*.
- **Idempotencia:** re-ejecuciones no duplican por dedup y *write mode* controlado.

**Nota.** El AIS puede contener omisiones o errores de origen; se prioriza evitar falsos positivos de integridad sin perder señales útiles para anomalías.



## 4. Proceso de desarrollo

### 4.1. Metodología

Esta sección cubre el **cómo** se ejecuta y opera el pipeline (orquestración, confiabilidad, despliegue y costos); las reglas de transformación y el esquema se tratan en las Secciones 3.3 y 1.

#### Proceso de desarrollo.

- **Ejecución y calendarización:** actualmente los disparadores son manuales o programados; aún no se emplea un flujo completamente *event-driven* (ver Sección 2.2).
- **Orquestración y repetibilidad:** los jobs se ejecutan en DataProc con tolerancia a fallos (speculation, retries) e *idempotencia* (sobrescritura controlada por partición y marcadores `_SUCCESS`). Esto permite relanzar sin duplicados ni corrupción de particiones.
- **Publicación a BigQuery:** una Cloud Function idempotente verifica/crea **datasets** y **tablas**; si no hay esquema, lo autodetecta a partir de una muestra en GCS. Su objetivo es que la capa *curated* quede lista para consulta por el dashboard.
- **Particionado y costos:** el diseño de partición por `ym`, permite escaneos selectivos y costos predecibles.
- **Versionamiento de dependencias:** se distribuyen *wheels* (p.ej., `pygeohash`) y un entorno reproducible (tarball del `venv`) a los ejecutores vía `spark.yarn.dist.archives`, asegurando UDFs consistentes.
- **Observabilidad:** se registran conteos globales y por día, además de tiempos por etapa, para detectar desviaciones (picos inusuales, lag de E/S). Los logs documentan fallas y se retienen para auditoría.

**Configuración del clúster (DataProc).** La capacidad de cómputo y la configuración actual del clúster es:

Cuadro 2: Tamaño del clúster de DataProc y características de nodos.

Rol	Tipo	vCPU	Memoria
Maestro	e2-highmem-4	4	32 GB
Trabajador	e2-highmem-2	2	16 GB
Trabajador	e2-highmem-2	2	16 GB
Trabajador	e2-highmem-2	2	16 GB
Trabajador	e2-highmem-2	2	16 GB

**Optimización y diferenciador.** La configuración *high-memory* actual mitiga la restricción operativa de procesar estrictamente mes a mes: ahora es posible ejecutar ventanas de *series de meses* más amplias con menor probabilidad de *spills*/OOM y menos corridas totales. No modifica la arquitectura, pero sí incrementa la capacidad efectiva del pipeline y mejora el tiempo de pared global.

### 4.2. Pruebas

En esta sección se documenta la **validación técnica del pipeline** (calidad de datos y comportamiento del sistema). Los objetivos de prueba fueron: (i) verificar que la ingesta y transformación preservan conteos y trazabilidad; (ii) garantizar que las reglas de limpieza y normalización operan según lo diseñado; (iii) asegurar idempotencia y reanudación segura por partición; y (iv) dejar evidencia auditable vía *logs* y marcadores de éxito.

### Validación de ingesta (*raw* → Parquet).

- **Conteos por día y por archivo:** para cada lote descargado se compararon los registros CSV con los generados en Parquet (mismo día *ymd*). Si los conteos no coincidían, el job abortaba y registraba la discrepancia.
- **Trazabilidad:** se propagaron metadatos de origen (*\_source\_file*, *\_ingest\_ts*) y se validó su presencia en la salida *raw* para auditoría.
- **Detección de duplicados en origen:** se verificó tamaño y nombre de archivo antes de descargar (evitando descargas repetidas) y se consolidó en GCS con estructura temporal por *ym*.

### Revisión de *outliers* y reglas de limpieza.

- **Velocidad:** se acotó SOG al rango físico [0, 70] nudos; se midió el % de registros recortados.
- **Coordenadas:** recorte de LAT a [−90, 90] y envoltura de LON a [−180, 180]; descarte de coordenadas nulas tras la corrección. Se verificó redondeo a 5 decimales de forma consistente.
- **Rumbos/estados:** normalización módulo 360° de COG/Heading; Heading=511 marcado como nulo. Mapas de NavStatusInt/Name con códigos fuera de catálogo etiquetados como *Unknown/Not reported*.

### Consistencia temporal.

- **UTC canónico:** BaseDateTime se parseó a TIMESTAMP en UTC y, a partir de éste, se derivaron *date*, *hour*, *week*, *month*, *quarter* y *ym*.
- **Integridad por ventana:** para cortes mensuales, se revisó que no existan *leaks* de registros fuera del *ym* correspondiente y que los límites de mes sean consistentes.

### Idempotencia y reanudación segura.

- **Clave de deduplicación:** MMSI, BaseDateTime como clave compuesta para garantizar series temporales sin duplicados en la capa *curated*.
- **Marcadores de éxito por partición:** escritura de *\_markers/ym=.../\_SUCCESS* al final de cada partición procesada. Si el job se relanza, sólo procesa particiones sin marcador, evitando sobrescrituras de trabajo correcto.
- **Relanzamiento seguro:** pruebas de re-ejecución completa y parcial (fallo inducido a mitad de mes) confirmaron que no se generan duplicados ni particiones corruptas.

### Observabilidad y evidencias.

- **Logs de DataProc/Spark:** conteos por etapa, tiempos de ejecución y tamaños de salida por *ym*; *speculation* y *retries* habilitados para robustez ante *stragglers*.
- **Métricas de control:** número de archivos procesados, % de registros filtrados por reglas (SOG fuera de rango, coordenadas inválidas), tasa de duplicados removidos y tiempo total por corrida.

Cuadro 3: Resumen de verificaciones de calidad del pipeline.

Chequeo	Criterio de aceptación	Evidencia
Conteos raw vs Parquet	Diferencia absoluta = 0 por <i>ymd</i>	Log de etapa e informe de conteos
Rango de SOG	[0, 70] nudos tras limpieza	Percentiles antes/después
Coordenadas válidas	LAT/LON dentro de dominio, sin nulos	Reporte de recortes/descartes
UTC y derivadas temporales	<i>hour</i> en 0 . . . 23; <i>ym</i> correcto	
Idempotencia por partición	Re-ejecución sin duplicados	Presencia de <i>_SUCCESS</i> y <i>diffs</i> nulos

**Validación previa en BigQuery.** Antes de integrar los resultados en la aplicación, se realizaron pruebas exploratorias en BigQuery para comprobar la validez de los *queries*, el esquema y la consistencia de la información retornada (rangos temporales, conteos y estadísticos básicos). Posteriormente, en la etapa de entrenamiento y evaluación de modelos de ML, los *queries* involucrados y los resultados producidos fueron verificados de igual manera, garantizando la funcionalidad para integrarlo en el dashboard. Como referencia visual, véanse las Figuras 2 y 3.

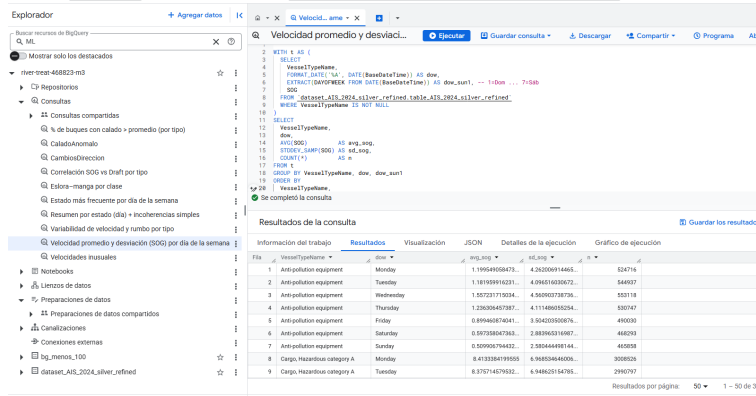


Figura 2: Pruebas y validaciones de *queries* en BigQuery (comprobación de periodos, filtros y métricas).

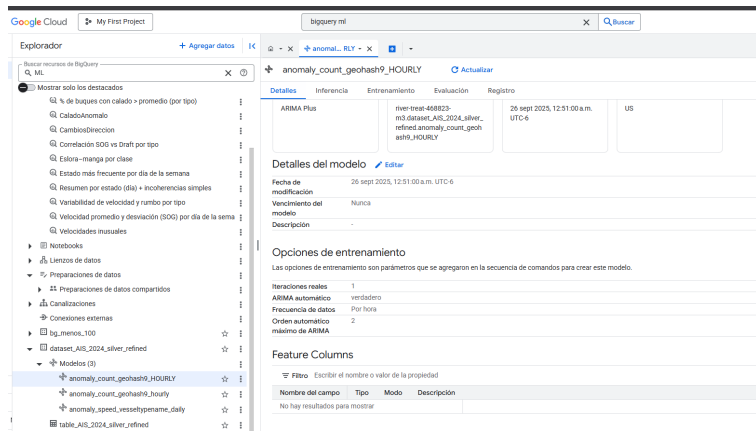


Figura 3: Verificación de entrenamiento/evaluación de modelos en BigQuery y revisión de resultados.

## 5. Resultados y conclusiones

Se desarrolló un *dashboard* interactivo que ejecuta las consultas validadas sobre los datos AIS y expone una superficie de entrenamiento y predicción basada en BigQuery ML (ARIMA PLUS). El tablero permite explorar coherencias/incoherencias de los mensajes, correlaciones geométricas y de desempeño, así como detectar comportamientos anómalos por serie y periodo. En conjunto, las visualizaciones y la interfaz de ML habilitan un flujo completo: (i) análisis exploratorio y verificación de datos; (ii) entrenamiento del modelo; (iii) evaluación de anomalías y revisión de resultados.

**Exploración y métricas.** Los módulos de exploración incluyen: (a) correlación *Eslora–Manga* por clase de buque, útil para detectar casos geométricamente atípicos; (b) *Calado anómalo* a partir de *z-score* por embarcación; (c) *Cambios de dirección* por encima de un umbral angular configurable, con soporte de filtro geográfico; y (d) variabilidad de velocidad y rumbo, más vistas complementarias (velocidad por día de semana, estados de navegación frecuentes). Estas vistas, además de presentar tablas paginadas, ofrecen gráficas de barras y de dispersión que facilitan identificar extremos y patrones.

**Entrenamiento y detección de anomalías.** La ventana de ML permite seleccionar métrica (*count* o *speed*), frecuencia (HOURLY/DAILY) e identificador de serie (*geohash9*, *VesselTypeName* o *MMSI*), ejecutar el entrenamiento y, posteriormente, lanzar la detección de anomalías sobre el rango de evaluación. Los resultados se listan en tabla y se visualizan con gráficos (barras por serie y dispersión temporal), con control del umbral de probabilidad para ajustar sensibilidad/especificidad. Este flujo fue verificado con consultas de prueba previas, asegurando que los *queries* de entrenamiento y predicción devuelven esquemas y valores consistentes antes de su integración final.

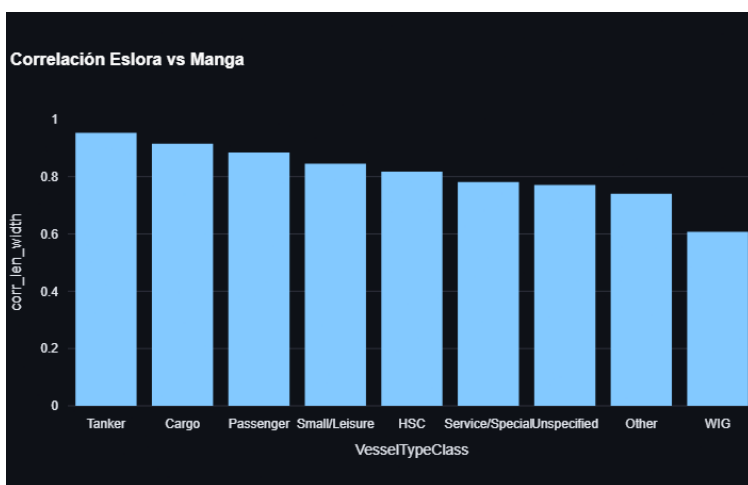


Figura 4: Correlación Eslora–Manga por clase de buque.

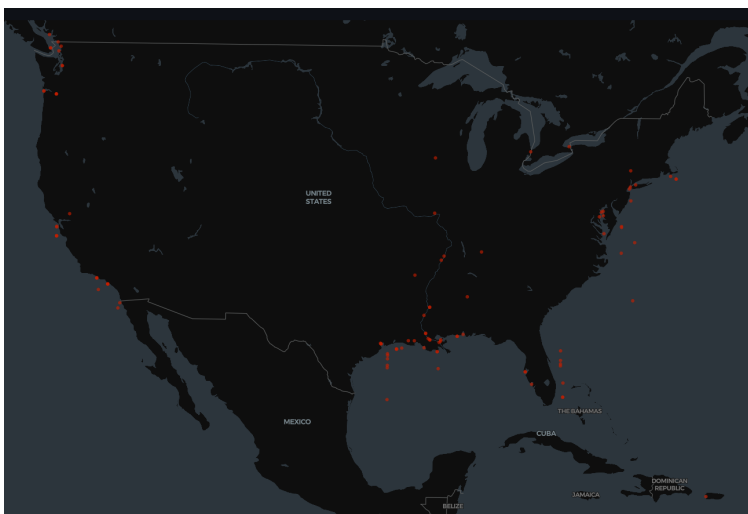


Figura 5: Eventos de cambio de dirección por encima de un umbral angular.

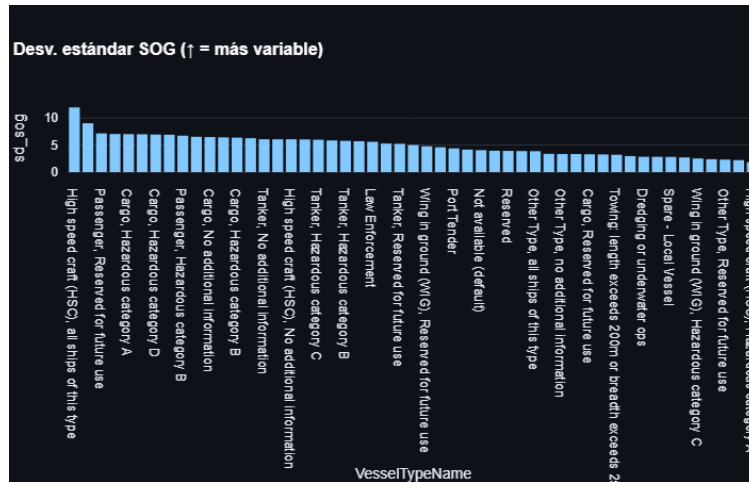


Figura 6: Variabilidad de velocidad y rumbo por tipo de buque.

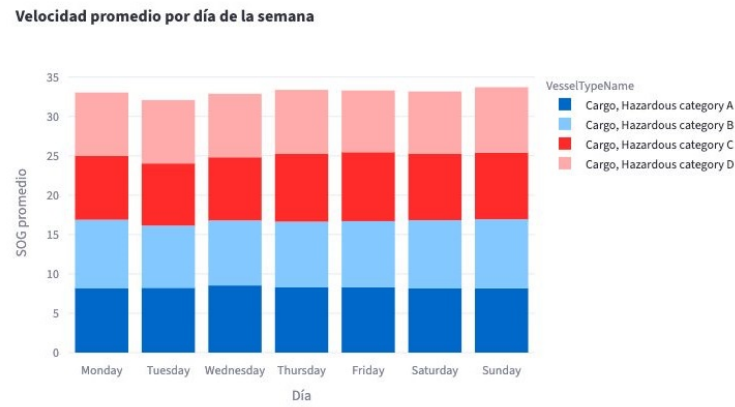


Figura 7: Velocidad promedio por día de la semana (barras apiladas por tipo de buque).

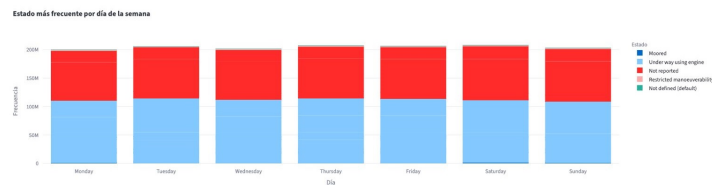


Figura 8: Estado de navegación más frecuente por día de la semana.

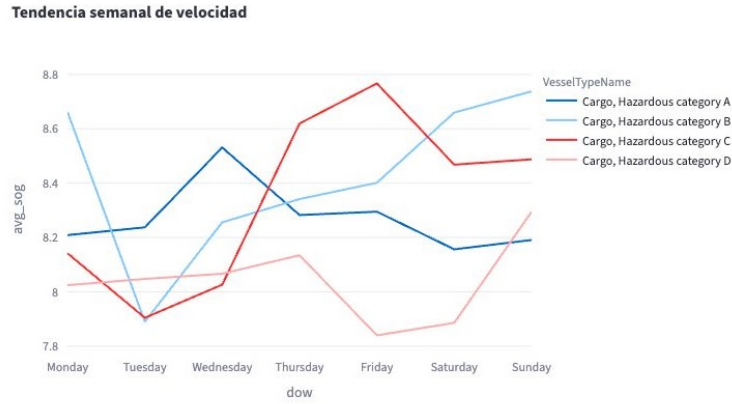


Figura 9: Tendencia semanal de velocidad por tipo de buque.



Figura 10: Ventana de entrenamiento y detección de anomalías con BigQuery ML.

## 5.1. Conclusión general

Se implementó un flujo *end-to-end* en Google Cloud que transforma datos masivos AIS en análisis operables y detección de anomalías: ingesta y preparación (GCE → GCS → DataProc/Spark), explotación analítica y de ML en BigQuery, y un *dashboard* en Streamlit con ventana de entrenamiento y predicción mediante BigQuery ML (ARIMA\_PLUS).

### Aportaciones principales.

- **Pipeline reproducible y escalable** con capas *raw/curated*, particionado temporal e idempotencia, optimizado para costos y tiempos de consulta.
- **Analítica descriptiva accionable** (patrones semanales de velocidad y estados, coherencias geométricas y de desempeño) para priorizar zonas, buques y ventanas temporales.
- **Detección de anomalías temporal** basada en series (ARIMA\_PLUS) con umbral configurable e integración directa en la interfaz para análisis y depuración de resultados.

### Lecciones y limitaciones.

- La **calidad y cobertura** del AIS condicionan los hallazgos; se requieren normalizaciones y validaciones previas en SQL.
- La versión actual opera en *batch* con disparadores manuales/programados y no incluye tiempo real.

- **Control de costos (GCS/BigQuery):** el volumen y la granularidad elevan costos de almacenamiento y consulta; se mitigarían con presupuestos/alertas de gasto.

## 6. Código utilizado

### 6.1. Link al repositorio con código fuente y salidas correspondientes

<https://github.com/enriuegomeztagle/MCD-BigData-SmartMaritimeTrafficMonitoring-FinalProject>