

## 1. Las 5 V's del Big Data:

**Volumen:** Amazon maneja cantidades enormes de datos todos los días. Esto incluye compras, búsquedas, reseñas, productos que los usuarios ven y también información logística. La escala es tan grande que llega a petabytes, y toda esa información es la base para personalizar la experiencia y optimizar los procesos.

**Velocidad:** Los datos en Amazon se procesan tanto en tiempo real como por lotes.

- **En tiempo real (Streaming):** Amazon procesa información de manera instantánea para mostrar recomendaciones mientras el usuario navega y para ajustar inventarios según la demanda del momento. También lo usan en logística para optimizar rutas y tiempos de entrega, reaccionando en segundos a cambios en la demanda o disponibilidad de productos.
- **Por lotes (Batch):** Además del tiempo real, también realizan procesamiento por lotes. Esto se aplica para analizar datos históricos y planificar la cadena de suministro, por ejemplo revisando grandes volúmenes de transacciones diarias o mensuales. Este análisis ayuda a detectar tendencias y preparar estrategias a largo plazo.

**Variedad:** Amazon utiliza una amplia variedad de tipos de datos, incluyendo:

- Estructurados: Incluyen datos como transacciones, inventario y precios. Son datos organizados en tablas que se pueden consultar y analizar fácilmente.
- No estructurados: Contienen información como reseñas, búsquedas, imágenes o videos. No tienen un formato fijo y requieren técnicas más avanzadas para analizarlos.
- Semi-estructurados: Son datos que tienen cierta organización pero no están en tablas fijas, como los registros de sensores en centros de distribución, datos de clics o logs de navegación.

**Veracidad:**

- Calidad de datos: Es fundamental que la información esté limpia, actualizada y sin errores. Con tantas fuentes distintas, un dato mal cargado o una reseña falsa puede impactar negativamente en las recomendaciones o en decisiones estratégicas.
- Confiabilidad: No toda la información se puede dar por cierta desde el inicio. Amazon necesita validar que los datos que utiliza realmente representen la realidad del negocio y el comportamiento de los usuarios para que las conclusiones sean correctas.

**Valor:** El Big Data le da a Amazon un valor enorme. Permite dar recomendaciones más precisas, optimizar la logística reduciendo tiempos y costos, y aumentar las ventas. Es uno de los pilares que les da una clara ventaja frente a sus competidores.

## 2. Almacenamiento:

Amazon utiliza una combinación de soluciones para poder manejar el enorme volumen y la variedad de datos que procesa a diario. Cada tipo de almacenamiento cumple un rol específico según la estructura de la información y la velocidad con la que se necesita acceder a ella.

- **Data Lake:** Guarda datos en su formato original y sin procesar, como información de sensores, clics o reseñas. Esto permite un esquema flexible y la posibilidad de analizarlos más adelante según sea necesario.
- **Data Warehouse:** Se utiliza para datos estructurados y ya procesados, como el historial de compras o las transacciones. Es la mejor opción cuando se necesitan análisis detallados o reportes.
- **Bases de datos NoSQL:** Manejan grandes volúmenes de datos no estructurados o semi-estructurados a alta velocidad. Amazon utiliza DynamoDB, inspirada en Google BigTable, que es ideal para este tipo de operaciones.

- **Bases de datos en memoria:** Procesan los datos directamente en RAM, lo que permite obtener resultados de forma inmediata. Esto es clave para las recomendaciones y análisis en tiempo real.
- **Sistemas de archivos distribuidos (HDFS):** Dividen y procesan datos en paralelo en múltiples servidores. Amazon lo implementa a través de Amazon Elastic MapReduce (EMR).

En cuanto a los desafíos, uno de los más grandes es la escalabilidad, ya que el volumen de datos crece rápidamente y los sistemas tradicionales no son suficientes. También está el costo, que aunque ha bajado gracias a la nube y las soluciones open source, puede aumentar si el almacenamiento no se gestiona correctamente y crece sin control.

### 3. Procesamiento y Análisis:

#### ¿Qué tipo de procesamiento se necesita (por lotes o en streaming)?

- Para las recomendaciones de productos y la detección de fraude, se necesita procesamiento en tiempo real (streaming), ya que la inmediatez de la respuesta es clave para la experiencia del usuario y la prevención de pérdidas.
- Para la gestión del inventario y la predicción de la demanda, si bien se utilizan datos en tiempo real de sensores, también se requieren procesos por lotes para el análisis de grandes volúmenes de datos históricos que no necesitan una respuesta inmediata, como las transacciones diarias procesadas al final del día para generar informes. La integración de ambos enfoques es fundamental.

#### ¿Qué herramientas de análisis serían las más adecuadas?

- **Aprendizaje Automático:** Base de los sistemas de recomendación de Amazon. Permite que los modelos aprendan de los datos para reconocer patrones y hacer predicciones personalizadas.
- **Minería de Datos:** Se utiliza para encontrar patrones y relaciones dentro de grandes volúmenes de información, lo que ayuda a entender mejor a los clientes y optimizar procesos.
- **Análisis Estadístico:** Sirve para validar resultados, segmentar clientes y construir modelos predictivos. Herramientas como R son muy usadas para trabajar con grandes conjuntos de datos.
- **Hadoop y MapReduce:** Frameworks que permiten procesar grandes volúmenes de datos de forma distribuida. Amazon lo implementa con Amazon EMR, dividiendo tareas complejas en subprocesos que se ejecutan en paralelo.
- **Bases de Datos NoSQL y en memoria:** Ideales para manejar grandes volúmenes de datos no estructurados y responder de forma rápida, especialmente en procesos que requieren inmediatez.
- **Lenguajes de programación:** Python, Scala y Java se utilizan para procesar y transformar datos, además de trabajar con plataformas como Apache Spark para análisis a gran escala.

### 4. Gobernanza y Seguridad:

¿Qué datos sensibles o personales podrían estar manejando? Amazon maneja una cantidad considerable de datos sensibles y personales de sus clientes, incluyendo:

- Historial de compras, búsquedas, productos vistos, reseñas leídas, artículos en el carrito no comprados.
- Preferencias y rechazos de perfiles de clientes en medios sociales.
- Datos de posición geográfica (geolocalización) y tiempo, si Amazon utiliza estos para servicios de entrega o análisis de comportamiento.

- Información de dispositivos (tipo de dispositivo, sistema operativo, etc.) para optimizar la experiencia móvil.
- Datos generados por sensores en sus operaciones.

**¿Qué desafíos de seguridad y privacidad tendrían que considerar para proteger la información?** La gestión de Big Data plantea muchas interrogantes en cuestiones de seguridad, privacidad y temas relacionados.

- **Privacidad del cliente:** Es fundamental que los datos personales no se utilicen ni se registren sin el consentimiento del usuario. También se deben respetar derechos como acceder a la información propia, corregirla o eliminarla. Las iniciativas de datos abiertos (Open Data) siempre deben proteger la privacidad.
- **Protección contra acceso no autorizado:** Se debe evitar que personas no autorizadas accedan a datos sensibles, sin importar dónde estén almacenados. Esto incluye enmascarar datos críticos y monitorear la actividad de las bases para detectar accesos sospechosos.
- **Cumplimiento normativo:** Amazon debe ajustarse a las leyes de protección de datos de cada país donde opera, como el GDPR en Europa y regulaciones locales en otros mercados.
- **Seguridad en entornos distribuidos y en la nube:** La expansión del Big Data y la nube requiere soluciones de seguridad específicas, ya que los datos pueden estar en servidores geográficamente dispersos y fuera del control físico directo de la empresa.
- **Riesgo de filtraciones:** El uso de dispositivos personales (BYOD) o aplicaciones no corporativas puede generar fugas de datos corporativos confidenciales si no se controla adecuadamente.
- **Calidad de los datos:** Aunque no es un problema de seguridad en sí, una mala calidad de datos puede llevar a decisiones equivocadas y sistemas poco confiables, afectando la eficacia de los procesos.
- **Gobierno de datos:** Debe formar parte de un marco más amplio de gobernanza, con normas, roles y responsabilidades claras, para optimizar y proteger la información durante todo su ciclo de vida.