

ANÁLISIS DE DATOS MASIVOS

MODELOS DE FLUJOS DE DATOS

Blanca Vázquez

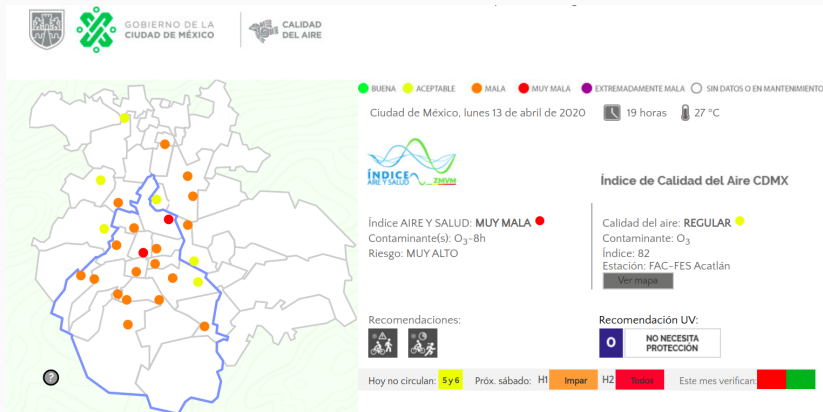
2 de octubre de 2024



Los sensores industriales pueden capturar grandes cantidades de datos

Imagen tomada de commons.wikimedia.org

ESTACIONES DE MONITOREO DE LA CALIDAD EL AIRE



ESTACIONES DE MONITOREO DE LA CALIDAD EL AIRE

O ₃																													
Fecha de consulta: 2020-04-13																													
Unidad de Medida: Partes por Millón (ppm)																													
Hora	A/M	ATI	BJU	CAM	CCA	CHO	CUA	CUT	FAC	FAR	GAM	HCM	IZT	LLA	LPR	MER	MGH	NEZ	PED	SAC	SAG	SFE	TAH	TLA	TLI	UIZ	UAX	VIF	XAL
1																													
2																													
3																													
4	0.043	0.017	0.014	0.010	0.018			0.002	0.019	0.026	0.016		0.013	0.005	0.014	0.004	0.020	0.023	0.025	0.033	0.024	0.039	0.020	0.016		0.016	0.023	0.004	
5	0.042	0.007	0.012	0.005	0.018			0.002	0.014	0.021	0.005		0.010	0.003	0.010	0.005	0.016	0.018	0.022	0.018	0.014	0.035	0.009	0.008		0.005		0.003	
6	0.038	0.001	0.008	0.005	0.007			0.002	0.008	0.011	0.004		0.003	0.003	0.001	0.004	0.004	0.014	0.010	0.007	0.007	0.022	0.016	0.002		0.004	0.012	0.002	
7	0.033	0.000	0.005	0.005	0.004	0.004	0.035	0.002	0.003	0.009	0.004	0.008	0.002	0.003	0.004	0.003	0.002	0.005	0.010	0.003	0.005	0.004	0.011	0.003	0.005	0.007	0.003		
8	0.035	0.001	0.006	0.007	0.011		0.034	0.002	0.005	0.011	0.013	0.014	0.008	0.004	0.007	0.008	0.008	0.004	0.016	0.005	0.006	0.007	0.011	0.004	0.009	0.008	0.010		
9	0.041	0.020	0.021	0.021	0.022		0.040	0.010	0.010	0.028	0.027	0.022	0.019	0.013	0.017	0.016	0.021	0.015	0.020	0.014	0.016	0.034	0.026	0.012		0.016	0.027	0.022	
10	0.049	0.031	0.038	0.034	0.050		0.042	0.035	0.022	0.042	0.043	0.035	0.037	0.035	0.022	0.036	0.034	0.034	0.046	0.035	0.027	0.044	0.050	0.028		0.037	0.043	0.044	
11	0.064	0.039	0.052	0.049	0.064		0.054	0.050	0.038	0.058	0.063	0.060	0.050	0.040	0.036	0.060	0.053	0.051	0.060	0.060	0.047	0.056	0.066	0.037		0.060	0.058	0.057	
12	0.079	0.046	0.061	0.062	0.077		0.070	0.058	0.047	0.068	0.076	0.073	0.062	0.048	0.056	0.075	0.069	0.057	0.081	0.066	0.050	0.072	0.077	0.051		0.066	0.068	0.057	
13	0.093	0.052	0.080	0.081	0.093		0.090	0.060	0.066	0.081	0.095	0.089	0.080	0.066	0.076	0.093	0.083	0.063	0.094	0.072	0.065	0.092	0.086	0.059		0.072	0.073	0.055	
14	0.097	0.060	0.086	0.097	0.100		0.099	0.062	0.078	0.093	0.107	0.110	0.091	0.082	0.064	0.103	0.107	0.083	0.103	0.095	0.081	0.118	0.077	0.078		0.095	0.089	0.065	
15	0.100	0.069	0.090	0.116	0.107		0.057	0.069	0.092	0.104	0.114	0.112	0.101	0.088	0.080	0.112	0.116	0.090	0.106	0.074	0.094	0.098	0.053	0.086		0.107	0.099	0.074	
16		0.072	0.094	0.121	0.105		0.073	0.085	0.073	0.111	0.117	0.113	0.100	0.098	0.094	0.112	0.110	0.063	0.100	0.063		0.072	0.054	0.091		0.087	0.096	0.087	
17	0.076	0.068	0.079	0.116	0.094		0.073	0.086	0.071	0.077	0.106	0.101	0.075	0.098	0.086	0.099	0.085	0.060	0.077	0.062	0.093	0.073	0.068	0.092		0.074	0.071	0.092	
18	0.069	0.059	0.062	0.082	0.073		0.061	0.062	0.086	0.059	0.067		0.064	0.062	0.056	0.071	0.069	0.055	0.065	0.063	0.052	0.062	0.079	0.077		0.071	0.070	0.080	
19	0.065	0.036	0.051	0.065	0.065		0.051	0.052	0.050	0.064	0.059		0.060	0.051	0.050	0.063	0.049	0.057	0.057	0.068	0.055	0.056	0.072	0.052		0.065	0.065	0.052	

Consulta el índice por zonas

Interpretación del Índice AIRE Y SALUD	
Concentraciones	Condición
0-0.051	Buena
>0.051 y 0.095	Aceptable
>0.095 y 0.135	Mala
>0.135 y 0.175	Muy Mala
>0.175	Extremadamente Mala
M	Mantenimiento

Los datos presentados en esta sección son preliminares y podrían sufrir modificaciones durante las siguientes etapas de validación.

Son datos que se generan constantemente (tiempo real) a partir de miles de fuentes de datos

- Normalmente los datos son enviados simultáneamente en conjuntos de tamaño pequeño (*kbs*)
- Si los datos no se almacenan o se procesan rápido, estos se perderán
- Dinámicos

- Atributos: cada atributo representa un tipo de dato (segmento, geo-localización, ID, ...)
- Marca de tiempo: indica hora y fecha de los datos generados
- Dato crudo: contiene la información original generada por la fuente de datos

- Monitoreo
- Dispositivos IoT
- Internet y tráfico web (por ej. secuencias de páginas visitadas (*clickstream*))
- Transacciones financieras
- Video juegos en línea
- Videos



Un sensor en el océano envía cada hora la temperatura del agua a una estación hidrológica (tasa de envío baja 4kb)

Un problema interesante sería:

- Un millón de sensores
- Cada uno enviando sus datos en una tasa de 10kb/segs
- Esto replicado cada 150 millas
- El océano Pacífico tiene 9,320.6 millas de norte a sur!!

DATOS DE IMÁGENES



Aproximadamente existen en órbita **5,000 satélites** que captan imágenes multispectrales de la Tierra de **resolución media y alta**.

Aproximadamente capturan y envían: **millón y medio de imágenes diarias**



Cámaras de vigilancia generalmente producen imágenes de baja resolución (en comparación con los satélites), sin embargo el intervalo de envío es de 1 segundo.

Londres tiene alrededor de **6 millones de cámaras**.

- Google procesa 81,226 búsquedas por segundo
- 3.5 miles de millones de búsquedas por día
- En 1999, a Google le tomó un mes indexar ≈ 50 millones de páginas. En el 2012 le tomó un minuto.
- Cada pregunta viaja 1,500 millas (hacia el centro de datos y de regreso)
- La respuesta a una consulta tarda 2 segundos (usando 1,000 computadoras)

¹Fuente: <https://www.internetlivestats.com/google-search-statistics/>

- Aplicaciones sencillas
 - Implementación de mínimo - máximo
 - Generación de informes básicos
 - Emitir alertas
- Aplicaciones complejas
 - Uso de aprendizaje máquina
 - Procesamiento de eventos y transmisiones

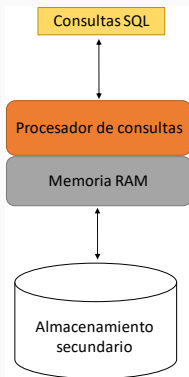
1. Memoria limitada para almacenar los datos
2. Debido a la vasta cantidad de datos, no es siempre posible generar respuestas exactas
3. Se espera que la calidad de la respuesta sea confiable
4. ¿Cómo trabajar con los datos (selección aleatoria, los últimos...)?

¿CÓMO SE PROCESAN LOS FLUJOS DE DATOS?

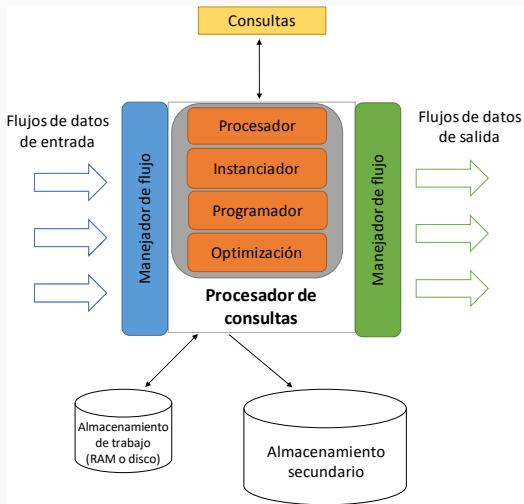
Un procesador de flujos de datos es un tipo de Sistema de Administración de datos (DSMS).

- Cualquier número de flujos puede ingresar al DSMS.
- Los flujos que se reciben no necesariamente deben tener la misma tasa de datos o tipo de datos
- El tiempo entre flujos no necesita ser uniforme.
- Los algoritmos para procesar los flujos puede involucrar resumen, filtrado o uso de ventanas.

MODELO GENERAL DE UN DBMS



MODELO GENERAL DE PROCESAMIENTO DE FLUJOS DE DATOS



- Las consultas son frecuentes
 - Los flujos son evaluados a medida que se van recibiendo
 - Actualizaciones constantes
- Las consultas son complejas
 - Pre-procesamiento de atributos y extracción de datos crudos

Existen dos formas generales para hacer consultas sobre los flujos de datos:

- Consultas permanentes: están almacenadas dentro del procesador, son ejecutadas permanentemente y producen salidas en momentos apropiados
- Consultas Ad-hoc: se realiza una sola vez sobre el flujo o flujos actuales

- Supongamos un sensor de temperatura en el océano, la consulta permanente sería “si la temperatura excede los 25 grados, emite una alerta”.
- Esta consulta solo dependen del último flujo recibido

- Otro ejemplo de consulta permanente sería: cada vez que llegue una nueva lectura (temp) genera el promedio de las últimas 24 lecturas
- Aquí almacenados las últimas 24 lecturas, cuando un nuevo valor llega se hace el cálculo y se borra la primera lectura

- Otro ejemplo de consulta permanente sería: obtén la temperatura máxima
- ¿Cómo lo haríamos?
- Y si la consulta es obtener el promedio, ¿cómo lo haríamos?

- Son consultas hechas una sola vez sobre los flujos actuales.
- Un enfoque común es almacenar una *ventana deslizando* de cada flujo en el *working storage*.

Técnica para el procesamiento de flujos de datos el cual divide dicho flujo en grupos de datos basándose en 2 parámetros

- Longitud de la ventana (*window length*): indica el tiempo que se tendrá en cuenta para el cálculo (desde t_{actual} hasta $t_{actual} - \text{longitud de ventana}$)
- Intervalo (*sliding interval*): cada cuánto tiempo se vuelve hacer los cálculos sobre los datos de la ventana

VENTANAS DESLIZANTES

Ejemplo: Actualizar cada segundo (intervalo) con el valor de la mayor compra de los últimos 2 segundos (longitud de la ventana)

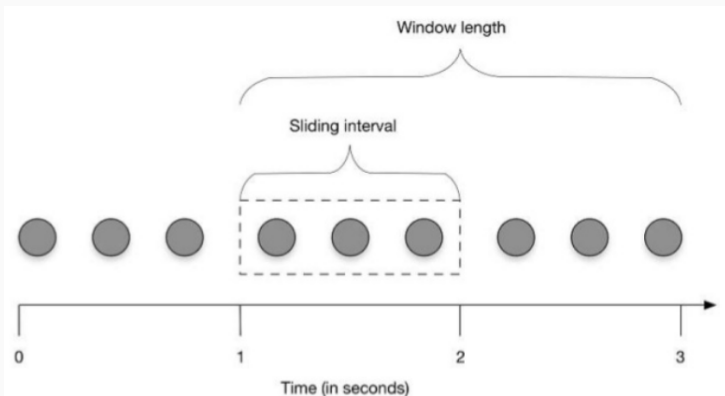
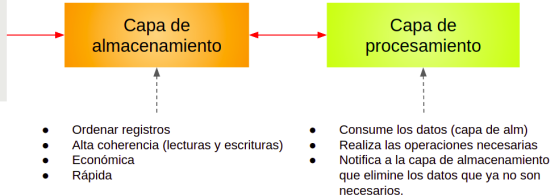
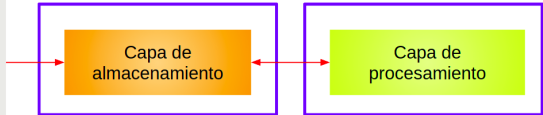


Imagen tomada de Workshop Apache Flink, 2016

CAPAS EN EL PROCESAMIENTO DEL FLUJO DE DATOS



CAPAS EN EL PROCESAMIENTO DEL FLUJO DE DATOS



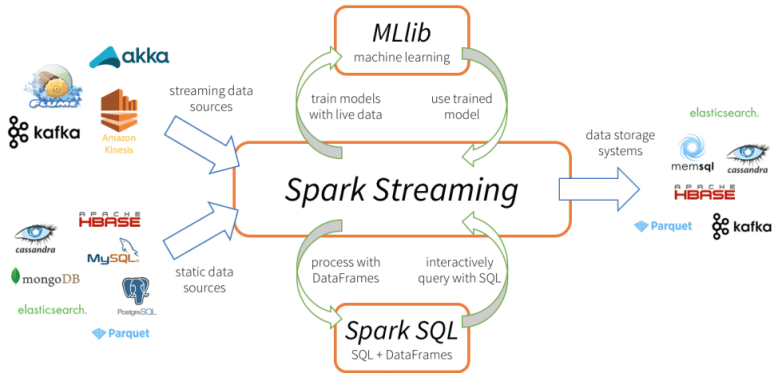
Una capa adicional se añade para cada una de las capas:

- Planificar la escalabilidad
- Durabilidad de los datos
- Tolerancia a fallos

Actualmente existen numerosas plataformas que soportan el procesamiento de flujos de datos

- Amazon Kinesis Streams
- Amazon Kinesis Firehose
- Apache Kafka
- Apache Flume
- Apache Spark Streaming
- Apache Storm

SPARK STREAMING



En Aprendizaje máquina ha surgido el *aprendizaje en línea*

- Nos permite modelar problemas en donde la entrada son flujos continuos de datos
- Se busca encontrar un algoritmo que aprenda a partir de los datos y que pueda adaptarse a pequeños cambios
- Ejemplos: Descenso del gradiente estocástico (SGD) permite pequeñas actualizaciones

COMPARACIÓN DBMS vs DSMS

DBMS	DSMS
Almacenamiento persistente	Almacenamiento transitorio
Acceso aleatorio	Acceso secuencial
Baja tasa de actualización	Tasas de múltiples Gbs
Servicios no de tiempo real	Servicios de tiempo real
Almacenamiento en disco ilimitada*	Memoria principal limitada