



Universidad de Valladolid

E.T.S Ingeniería Informática

Trabajo Fin de Grado

Grado en Ingeniería Informática,
mención en Computación

Algoritmos para Big Data

Autor:
Sergio García Prado



Universidad de Valladolid

E.T.S Ingeniería Informática

Trabajo Fin de Grado

Grado en Ingeniería Informática,
mención en Computación

Algoritmos para Big Data

Autor:

Sergio García Prado

Tutor:

Manuel Barrio Solórzano

Prefacio

Para entender el contenido de este documento así como la metodología seguida para su elaboración, se han de tener en cuenta diversos factores, entre los que se encuentran el contexto académico en que ha sido redactado, así como el tecnológico y social. Es por ello que a continuación se expone una breve descripción acerca de los mismo, para tratar de facilitar la comprensión sobre el alcance de este texto.

Lo primero que se debe tener en cuenta es el contexto académico en que se ha llevado a cabo. Este documento se ha redactado para la asignatura de **Trabajo de Fin de Grado (mención en Computación)** para el *Grado de Ingeniería Informática*, impartido en la *E.T.S de Ingeniería Informática* de la *Universidad de Valladolid*. Dicha asignatura se caracteriza por ser necesaria la superación del resto de las asignaturas que componen los estudios del grado para su evaluación. Su carga de trabajo es de **12 créditos ECTS**, cuyo equivalente temporal es de *300 horas* de trabajo del alumno, que se han llevado a cabo en un periodo de 4 meses.

La temática escogida para realizar dicho trabajo es **Algoritmos para Big Data**. El Big Data es la disciplina que se encarga de “todas las actividades relacionadas con los sistemas que manipulan grandes conjuntos de datos. Las dificultades más habituales vinculadas a la gestión de estas cantidades de datos se centran en la recolección y el almacenamiento, búsqueda, compartición, análisis, y visualización. La tendencia a manipular enormes cantidades de datos se debe a la necesidad en muchos casos de incluir dicha información para la creación de informes estadísticos y modelos predictivos utilizados en diversas materias.” [Wik17]

Uno de los puntos más importantes para entender la motivación por la cual se ha escogido dicha temática es el contexto tecnológico en que nos encontramos. Debido a la importante evolución que están sufriendo otras disciplinas dentro del mundo de la informática y las nuevas tecnologías, cada vez es más sencillo y económico recoger gran cantidad de información de cualquier proceso que se dé en la vida real. Esto se debe a una gran cantidad de factores, entre los que se destacan los siguientes:

- **Reducción de costes derivados de la recolección de información:** Debido a la constante evolución tecnológica cada vez es más barato disponer de mecanismos (tanto a nivel de hardware como de software), a partir de los cuales se puede recabar datos sobre un determinado suceso.
- **Mayor capacidad de cómputo y almacenamiento:** La recolección y manipulación de grandes cantidades de datos que se recogen a partir de sensores u otros métodos requieren por tanto del apoyo de altas capacidades de cómputo y almacenamiento. Las tendencias actuales se están apoyando en técnicas de virtualización que permiten gestionar sistemas de gran tamaño ubicados en distintas zonas geográficas como una unidad, lo cual proporciona grandes ventajas en cuanto a reducción de complejidad algorítmica a nivel de aplicación.
- **Mejora de las telecomunicaciones:** Uno de los factores que ha permitido una gran disminución de la problemática relacionada con la virtualización y su capacidad de respuesta ha sido el gran avance a nivel global que han sufrido las telecomunicaciones en los últimos años, permitiendo disminuir las barreras geográficas entre sistemas tecnológicos dispersos.

Gracias a este conjunto de mejoras se ha llegado al punto en que existe la oportunidad de poder utilizar una gran cantidad de conocimiento, que individualmente o sin un apropiado procesamiento, carece de valor a nivel de información.

El tercer factor que es necesario tener en cuenta es la tendencia social actual, que cada vez más, está concienciada con el valor que tiene la información. Esto se ve reflejado en un amplio abanico de aspectos relacionados con el comportamiento de la población:

- **Monitorización de procesos laborales:** Muchas empresas están teniendo en cuenta la mejora de la productividad de sus empleados y máquinas. Por tanto, buscan nuevas técnicas que les permitan llevar a cabo dicha tarea. En los últimos años se ha dedicado mucho esfuerzo en implementar sistemas de monitorización que permitan obtener información para después procesarla y obtener resultados valiosos para dichas organizaciones.
- **Crecimiento exponencial de las plataformas de redes sociales:** La inherente naturaleza social del ser humano hace necesaria la expresión pública de sus sentimientos y acciones, lo cual, en el mundo de la tecnología se ha visto reflejado en un gran crecimiento de las plataformas de compartición de información así como de las de comunicación.
- **Iniciativas de datos abiertos por parte de las administraciones públicas:** Muchas insituciones públicas están dedicando grandes esfuerzos en hacer visible la información que poseen, lo que conlleva una mejora social aumentando el grado de transparencia de las mismas, así como el nivel de conocimiento colectivo, que puede ser beneficioso tanto para ciudadanos como para empresas.

Como consecuencia de este cambio social, posiblemente propiciado por el avance tecnológico anteriormente citado, la población tiene un mayor grado de curiosidad por aspectos que antes no tenía la capacidad de entender, debido al nivel de complejidad derivado del tamaño de los conjuntos de muestra necesarios para obtener resultados fiables.

En este documento no se pretenden abordar temas relacionados con las técnicas utilizadas para recabar nuevos datos a partir de los ya existentes. A pesar de ello se realizará una breve introducción sobre dicho conjunto de estrategias, entre las que se encuentran: *Heurísticas*, *Regresión Lineal*, *Árboles de decisión*, *Máquinas de Vector Soporte (SVM)* o *Redes Neuronales Artificiales*.

Por contra, se pretende realizar un análisis acerca de los diferentes algoritmos necesarios para manejar dichas cantidades ingentes de información, en especial de su manipulación a nivel de operaciones básicas, como operaciones aritméticas, búsqueda o tratamiento de campos ausentes. Para ello, se tratará de acometer dicha problemática teniendo en cuenta estrategias de paralelización, que permitan aprovechar en mayor medida las capacidades de cómputo existentes en la actualidad.

Otro de los aspectos importantes en que se quiere orientar este trabajo es el factor dinámico necesario para entender la información, lo cual conlleva la búsqueda de nuevas estrategias algorítmicas de procesamiento en tiempo real. Por lo tanto, se pretende ilustrar un análisis acerca de las soluciones existentes en cada caso con respecto a la solución estática indicando las ventajas e inconvenientes de la versión dinámica según corresponda.

Índice general

Prefacio	1
1. Introducción	5
2. Algoritmos para Streaming	7
2.1. Introducción	7
2.2. Modelo en Streaming	9
2.3. Estructura básica	12
2.4. Medidas de Análisis	13
2.5. Algoritmo de Morris	17
2.6. Algoritmo de Flajolet-Martin	17
2.7. Aproximación a los Momentos de Frecuencia	17
3. Estructuras de Datos de Resumen	19
3.1. Introducción	19
3.2. Tipos de Estructuras de Datos de Resumen	19
3.3. Sketching	19
3.4. Count-Min Sketch	19
3.5. Count Sketch	19
3.6. AMS Sketch	19
3.7. HyperLogLog	19
4. Algoritmos para Grafos	21
4.1. Introducción	21
4.2. Modelo en Semi-Streaming	21
5. Reducción de la Dimensionalidad	23
5.1. Introducción	23
5.2. Teorema de Johnson-Lindenstrauss	23
5.3. Búsqueda de Vecinos más Cercanos	23
6. Técnicas de Minería de Datos	25
6.1. Introducción	25
6.2. Aprendizaje Supervisado y No Supervisado	25
6.3. Árboles de Decisión	25
6.4. Regresión Lineal	25
6.5. Redes Neuronales	25
6.6. Manifold Learning	25

7. Paralelización a Gran Escala	27
7.1. Introducción	27
7.2. Sistemas de Ficheros Distribuidos	27
7.3. Modelo de acceso a Memoria	27
7.4. Complejidad de la Compucación	27
7.5. MapReduce	27
A. Metodología de Trabajo	29
B. ¿Cómo ha sido generado este documento?	31
Bibliografía	33

Capítulo 1

Introducción

[TODO]

Capítulo 2

Algoritmos para Streaming

2.1. Introducción

En este capítulo se trata de realizar una descripción en profundidad acerca de los *Algoritmos en Streaming* desde una perspectiva tanto a teórica como práctica. Para ello se describirá el modelo de cómputo en que se enmarcan dichos algoritmos (Modelo en Streaming) en la sección 2.2 además de su estructura básica en la sección 2.3. El motivo de dicha descripción se debe a que los *Algoritmos para Streaming* presentan un conjunto de peculiaridades respecto de la gran mayoría de algoritmos utilizados comunmente. [TODO introducir el resto de secciones]

Para realizar una primera aproximación acerca de en qué consiste esta categoría algorítmica es necesario realizar una diferenciación entre distintos conceptos relacionados con ella, que pueden producir confusiones debido a su similitud o abusos previos del lenguaje. Por lo tanto, a continuación se describen conceptos relacionados con los *Algoritmos en Streaming* que permitirán introducir al lector en el contexto del problema. Además, se realiza una diferenciación acerca de los factores que se pretenden optimizar a partir de esta estrategia de diseño de algoritmos.

2.1.1. Computación en Tiempo Real

El primer concepto que se describe es **Computación en Tiempo Real**, que tal y cómo describen Shin y Ramanathan [SR94] se caracteriza por tres términos que se describen a continuación:

- **Tiempo**(*time*): En la disciplina de *Computación en Tiempo Real* el tiempo de ejecución de una determinada tarea es especialmente crucial para garantizar el correcto desarrollo del cómputo, debido a que se asume un plazo de ejecución permitido, a partir del cual la solución del problema deja de tener un valor óptimo. Shin y Ramanathan[SR94] diferencian entre tres categorías dentro de dicha restricción, a las cuales denominan *hard*, *firm* y *soft*, dependiendo del grado de relajación de la misma.
- **Confiabilidad**(*correctness*): Otro de los puntos cruciales en un sistema de *Cómputación en Tiempo Real* es la determinación de una unidad de medida o indicador acerca de las garantías de una determinada solución algorítmica para cumplir lo que promete de manera correcta en el tiempo esperado.
- **Entorno**(*environment*): El último factor que indican Shin y Ramanathan[SR94] para describir un sistema de *Computación en Tiempo Real* es el entorno del mismo, debido a que este condiciona el conjunto de tareas y la periodicidad en que se deben llevar a cabo. Debido a esta razón, realizan una diferenciación entre: *a*) tareas periódicas *periodic tasks* las cuales se realizan secuencialmente a partir de la finalización de una ventana de tiempo, y *b*) tareas no periódicas *aperiodic tasks* que se llevan a cabo debido al suceso de un determinado evento externo.

2.1.2. Problemas Dinámicos

Una vez completada la descripción acerca de lo que se puede definir como *Computación en Tiempo Real*, conviene realizar una descripción desde el punto de vista de la *teoría de complejidad computacional*. Para definir este tipo de problemas, se utiliza el término *problemas dinámicos*, los cuales consisten en aquellos en los cuales es necesario recalcular su solución conforme el tiempo avanza debido a variaciones en los parámetros de entrada del problema (Nótese que dicho término no debe confundirse con la estrategia de *programación dinámica* para el diseño de algoritmos).

Existen distintas vertientes dependiendo del punto de vista desde el que se estudien, tanto de la naturaleza del problema (soluciones dependientes temporalmente unas de otras o soluciones aisladas) como de los parámetros de entrada (entrada completa en cada nueva ejecución o variación respecto de la anterior). Los *Algoritmos para Streaming* están diseñados para resolver *problemas dinámicos*, por lo que en la sección 2.2, se describe en profundidad el modelo en que se enmarcan.

A continuación se indican los principales indicadores utilizados para describir la complejidad de una determinada solución algorítmica destinada a resolver un problema de dicha naturaleza:

- **Espacio:** Cantidad de espacio utilizado en memoria durante la ejecución del algoritmo.
- **Inicialización:** Tiempo necesario para la inicialización del algoritmo.
- **Procesado:** Tiempo necesario para procesar una determinada entrada.
- **Pregunta[TODO Buscar mejor palabra]:** Tiempo necesario para procesar la solución a partir de los datos de entrada procesados hasta el momento.

2.1.3. Algoritmos Online vs Algoritmos Offline

Una vez descrita la problemática de *Computación en Tiempo Real* en la sección 2.1.1 y la categoría de *Problemas Dinámicos* en la sección 2.1.2, en esta sección se pretende ilustrar la diferencia entre los *Algoritmos Online* y los *Algoritmos Offline*. Para ello, se ha seguido la diferenciación propuesta por Karp [Kar92], en la cual se plantea el problema de la siguiente manera (Se utilizará la misma notación que sigue Muthukrishnan[Mut05] para tratar mantener la consistencia durante todo el documento): Sea A el conjunto de datos o eventos de entrada, siendo cada $A[i]$ el elemento i -ésimo del conjunto. En el caso de los *Algoritmos Online* supondremos que es el elemento recibido en el instante i . A continuación se muestran las características de cada subgrupo:

- **Algoritmos Offline:** Esta categoría contiene todos los algoritmos que realizan el cómputo suponiendo el acceso a cualquier elemento del conjunto de datos A durante cualquier momento de su ejecución. Además, en esta categoría se impone la restricción de que el A debe ser invariante respecto del tiempo, lo que impone que para la adaptación del resultado a cambios, este tenga que realizar una nueva ejecución desde su estado inicial. Nótese que por tanto, dentro de este grupo se engloba la mayoría de algoritmos utilizados comúnmente.
- **Algoritmos Online:** Son aquellos que calculan el resultado a partir de una secuencia de sucesos $A[i]$, los cuales generan un resultado dependiente del valor, y posiblemente de los sucedidos anteriormente. A partir de dicha estrategia, se añade una componente dinámica, la cual permite que tamaño del conjunto de datos de entrada A no tenga impuesta una restricción acerca de su longitud. Por contra, en este modelo no se permite conocer el suceso $A[i + 1]$ en el momento i . Esto encaja perfectamente en el modelo que se describirá en la sección 2.2.

Según la diferenciación que se acaba de describir, estas dos estrategias de diseño de algoritmos encajan en disciplinas distintas de diseño de algoritmos, teniendo una gran ventaja a nivel de eficiencia en el caso estático los *Algoritmos Offline*, pero quedando prácticamente inutilizables cuando la computación es en tiempo real, donde es mucho más apropiado el uso de estrategias de diseño de *Algoritmos Online*.

Como medida de eficiencia para los *Algoritmos Online*, Karp [Kar92] propone el **Ratio Competitivo**, el cual se define como la cota inferior del coste de cualquier nueva entrada con respecto de la que tiene menor coste. Sin embargo, dicha medida de eficiencia no es comúnmente utilizada en el caso de los *Algoritmos para Streaming* por la componente estocástica de los mismos, para los cuales son más apropiadas medidas probabilistas. A continuación se describen las ventajas de estos respecto de su vertiente estática.

2.1.4. Algoritmos Probabilistas

Los *Algoritmos Probabilistas* son una estrategia de diseño que emplea en un cierto grado de aleatoriedad en alguna parte de su lógica. Estos utilizan distribuciones uniformes de probabilidad para tratar de conseguir un incremento del rendimiento en su caso promedio. A continuación se describen los dos tipos de algoritmos probabilísticos según la clasificación realizada por Babai [Bab79]:

- **Algoritmos Las Vegas:** Devuelven un resultado incorrecto con una determinada probabilidad, pero avisan del resultado incorrecto cuando esto sucede. Para contrarrestar este suceso basta llevar a cabo una nueva ejecución del algoritmo, lo cual tras un número indeterminado de ejecuciones produce un resultado válido.
- **Algoritmos Monte Carlo:** Fallan con un cierto grado de probabilidad, pero en este caso no avisan del resultado incorrecto. Por lo tanto, lo único que se puede obtener al respecto es un indicador de la estimación del resultado correcto hacia la que converge tras varias ejecuciones. Además, se asegura una determinada cota del error ϵ , que se cumple con probabilidad δ .

La razón anecdótica por la cual Babai [Bab79] decidió denominar dichas categorías de algoritmos de esta manera se debe a lo siguiente (teniendo en cuenta el contexto de lengua inglesa): cuando se va a un casino en *Las Vegas* y se realiza una apuesta el *croupier* puede decir si se ha ganado o perdido porque habla el mismo idioma. Sin embargo, si sucede la misma situación en *Monte Carlo*, tan solo se puede conocer una medida de probabilidad debido a que en este caso el *croupier* no puede comunicarlo por la diferencia dialéctica.

2.1.5. Algoritmos Online Probabilistas vs Deterministas

La idea subyacente acerca del diseño de los *Algoritmos Online* es la mejora de eficiencia con respecto a sus homónimos estáticos cuando el conjunto de valores de entrada es dependiente de los resultados anteriores. Sin embargo, existen casos en que la frecuencia de ejecución del algoritmo, debido a una alta tasa de llegada de sucesos de entrada, las soluciones deterministas se convierten en alternativas poco escalables.

Dicha problemática se ha incrementado de manera exponencial debido al avance tecnológico y la gran cantidad de información que se está generando en la actualidad. Este fenómeno ha convertido en algo necesario el diseño de estrategias basadas en sucesos probabilísticos que reducen en gran medida el coste computacional eliminando el determinismo de la solución.

2.2. Modelo en Streaming

En esta sección se describen los aspectos formales del *Modelo en Streaming*. Para ello se ha seguido la representación definida por Muthukrishnan [Mut05]. Lo primero por tanto, es definir lo que es un flujo de datos o *Data Stream* como una “secuencia de señales digitalmente codificadas utilizadas para representar una transmisión de información” [Ins17]. Muthukrishnan [Mut05] hace una aclaración sobre dicha definición y añade la objeción de que los datos de entrada deben tener un ritmo elevado de llegada. Debido a esta razón existe complejidad a tres niveles:

- **Transmisión:** Ya que debido a la alta tasa de llegada es necesario diseñar un sistema de interconexiones que permita que no se produzcan congestiones debido a la obtención de los datos de entrada.
- **Computación:** Puesto que la tarea de procesar la gran cantidad de información que llega por unidad de tiempo produce cuellos de botella en el planteamiento. Por lo que es necesario implementar técnicas algorítmicas con un reducido nivel de complejidad computacional que para contrarrestar dicha problemática.
- **Almacenamiento:** Debido a la gran cantidad de datos que se presentan en la entrada, deben existir técnicas que permitan almacenar dicha información de manera eficiente. Esto puede ser visto desde dos puntos de vista diferentes: *a)* tanto desde el punto de vista del espacio, tratando de minimizar el tamaño de los datos almacenados, maximizando la cantidad de información que se puede recuperar de ellos, *b)* como desde el punto de vista del tiempo necesario para realizar operaciones de búsqueda, adición, eliminación o edición.

[TODO ejemplo de los tres niveles de complejidad a partir del análisis meteorológico]

2.2.1. Formalismo para Streaming

Una vez descritos los niveles de complejidad a los que es necesario hacer frente para abordar problemas en el *Modelo en Streaming*, se realiza una descripción de los distintos modelos que propone Muthukrishnan [Mut05] en los apartados 2.2.2, 2.2.3 y 2.2.4. La especificación descrita en dichos apartados será seguida durante el resto del capítulo. Para ello nos basaremos en el siguiente formalismo:

Sea $a_1, a_2, \dots, a_t, \dots$ un flujo de datos de entrada (*Input Stream*), de tal manera que cada elemento debe presentar un orden de llegada secuencial respecto de $t \in \mathbb{N}$. Esto también se puede ver de la siguiente manera: el elemento siguiente a la llegada de a_{t-1} debe ser a_t y, por inducción, el próximo será a_{t+1} . Es necesario aclarar que t no se refiere a unidades propiamente temporales, sino a la posición en la entrada.

$$\mathbf{A}_t : [1 \dots N] \rightarrow \mathbb{R}^2 \quad (2.1)$$

El siguiente paso para describir el formalismo es añadir la función \mathbf{A}_t , cuyo dominio e imagen se muestran en la ecuación (2.1). Esta función tiene distintas interpretaciones dependientes del *Modelo en Streaming* con el cual se esté tratando en cada caso, pero la idea subyacente se puede resumirse asumiendo que la primera componente almacena el valor, mientras que la segunda almacena el número de ocurrencias de dicho valor. Algo común a todos ellos es la variable t , que se corresponde con resultado de la función en el instante de tiempo t . Por motivos de claridad, en los casos en que nos estemos refiriendo un único momento, dicha variable será obviada en la notación.

2.2.2. Modelo de Serie Temporal

El *Modelo de Serie Temporal* o *Time Series Model* se refiere, tal y como indica su nombre, a una serie temporal, es decir, modeliza los valores que toma la variable i respecto de t , codificados en el modelo como $a_t = (i, 1)$. Nótese que se utiliza el valor 1 en la segunda componente de a_t , la razón de ello se debe a la definición de la imagen de \mathbf{A} en la ecuación (2.1). A pesar de ello, dicho campo es irrelevante en este modelo, por lo que se podría haber escogido cualquier otro arbitrariamente. La razón por la cual se ha utilizado el valor 1 ha sido el refuerzo de la idea de que en este caso, el valor que toma la a_t en un determinado momento, no volverá a variar su valor, pero quedará obsoleto con la llegada de a_{t+1} .

El modelo se describe de manera matemática mediante la función \mathbf{A} , tal y como se ilustra en la ecuación (2.2). Textualmente, esto puede traducirse diciendo que la función \mathbf{A} representa una estructura de datos que almacena el valor de todos los elementos recibidos en la entrada hasta el instante de tiempo t , es decir, actúa como un historial. Un ejemplo de este modelo son los valores en bolsa que toma una determinada empresa a lo largo del tiempo.

$$\mathbf{A}(t) = a_t \quad (2.2)$$

2.2.3. Modelo de Caja Registradora

El *Modelo de Caja Registradora* o *Cash Register Model* consiste en la recepción de incrementos de un determinado valor i . El nombre del modelo hace referencia al funcionamiento de una caja registradora (suponiendo que el pago se realiza de manera exacta), que recibe billetes o monedas de tipos diferentes de manera secuencial.

Para describir dicho modelo, previamente hay que realizar una aclaración acerca del contenido del elemento $a_t = (i, I_t)$, de manera que i representa el valor recibido, mientras que $I_t \geq 0$ indica el incremento en el instante t . Una vez aclarada esta definición, la función \mathbf{A}_t , se construye tal y como se indica en la ecuación (2.3).

$$\mathbf{A}_t(i) = A_{t-1}(i) + I_t \quad (2.3)$$

El *Modelo de Caja Registradora* es ampliamente utilizado en la formalización de problemas reales debido a que muchos fenómenos siguen esta estructura. Un ejemplo de ello es el conteo de accesos a un determinado sitio web, los cuales se corresponden con incrementos I_t , en este caso de carácter unitario realizados por un determinado usuario i en el momento t .

2.2.4. Modelo de Molinete

El *Modelo de Molinete* o *Turnstile Model* se corresponde con el caso más general, en el cual no solo se permiten incrementos, sino que también se pueden realizar decrementos en la cuenta. El nombre que se le dió a este modelo se debe al funcionamiento de los molinetes que hay en las estaciones de metro para permitir el paso a los usuarios, que en la entrada incrementan la cuenta del número de personas, mientras que en la salida los decrementan. La relajación originada por la capacidad de decremento ofrece una mayor versatilidad, que permite la contextualización de un gran número de problemas en este modelo. Por contra, añade un numerosas complicaciones a nivel computacional, tal y como se verá a lo largo del capítulo.

Tal y como ocurre en el caso anterior, para describir este modelo, lo primero es pensar en la estructura de los elementos en la entrada, que están formados por $a_t = (i, U_t)$, algo muy semejante a lo descrito en el *Modelo de Caja Registradora*. Sin embargo, en este caso U_t no tiene restricciones en su imagen, sino que puede tomar cualquier valor tanto positivo como negativo, lo cual añade el concepto de decremento. La construcción de la función \mathbf{A}_t se describe en la ecuación (2.4).

$$\mathbf{A}_t(i) = A_{t-1}(i) + U_t \quad (2.4)$$

Muthukrishnan [Mut05] hace una diferenciación dentro de este modelo dependiendo del nivel de exigencia que se le pide al modelo, se dice que es un *Modelo de Molinete estricto* cuando se añade la restricción $\forall i, \forall t \mathbf{A}_t(i) \geq 0$, mientras que se dice que es un *Modelo de Molinete relajado* cuando dicha restricción no se tiene en cuenta.

Un ejemplo de este modelo es el conteo del número de usuarios que están visitando un determinado sitio web, tomando U_t el valor 1 en el caso de una nueva conexión y -1 en el caso de una desconexión. En este ejemplo el valor i representa una determinada página dentro del sitio web.

2.3. Estructura básica

Puesto que la naturaleza intrínseca de los *Algoritmos para Streaming* hace que procesen los elementos de entrada según van llegando, esto presenta peculiaridades con respecto a otras categorías algorítmicas utilizadas comúnmente. Por tanto, primero se describirá la estructura básica que siguen los algoritmos más comúnmente utilizados para después mostrar la estrategia seguida en el caso de Streaming.

Los algoritmos clásicamente estudiados para resolver la mayoría de problemas se basan en la idea de funciones matemáticas. Es decir, se les presenta un conjunto de valores en la entrada, y a partir de ellos, realizan una determinada transformación sobre ellos, que genera como resultado una determinada salida. Nótese que esta idea no genera ninguna restricción acerca de lo que puede suceder en dicho proceso, es decir, no se restringe el uso de estructuras de datos auxiliares o técnicas similares.

Esta visión no se enmarca correctamente en el contexto de los *Algoritmos para Streaming*. La razón se debe a que la entrada no es propiamente un conjunto de datos, sino que se refiere a un flujo en sí mismo. Esta característica conlleva que en un gran número de ocasiones ya no sea necesario obtener los resultados tras cada llamada al algoritmo, ya que estos podrían carecer de interés o requerir un sobrecoste innecesario. Por lo tanto, el concepto de función matemática pierde el sentido en este caso, ya que estas exigen la existencia de un valor como resultado.

Un concepto más acertado para modelizar un *Algoritmo para Streaming* podría ser lo que en los lenguajes de programación derivados de *Fortran* se denomina subrutina, es decir, una secuencia de instrucciones que realizan una tarea encapsulada como una unidad. Sin embargo, para poder describir correctamente la estructura de un *Algoritmo para Streaming* hace falta algo más. La razón de ello es que a partir de dicho modelo de diseño no sería posible realizar peticiones acerca de lo calculado, es decir, sería una estrategia puramente procedural. Para corregir dicha problemática surge el concepto de *query* o pregunta. A través de dicha idea se pretende representar la manera de obtener un resultado a partir del cómputo realizado hasta el momento actual.

En resumen, con dicha estrategia de diseño se consigue separar la parte de procesamiento de la entrada con la parte de consulta del resultado, lo cual proporciona amplias ventajas para el modelo seguido por los *Algoritmos para Streaming*. Sin embargo, dicha estrategia produce un sobrecoste espacial con respecto del modelo de algoritmo clásico. Este se debe a la necesidad de mantener una estructura de datos en la cual se almacenen los resultados parciales referentes al flujo de entrada.

Los algoritmos *Algoritmos para Streaming* se componen por tanto de algoritmo de procesamiento del flujo de datos, una estructura de datos que almacena dichos resultados, y por último, un algoritmo de procesamiento de la *query* o pregunta necesaria para obtener los resultados requeridos. a continuación se dividen las fases para el funcionamiento de un algoritmo de dichas características.

- **Inicialización:** En esta fase se llevan a cabo el conjunto de tareas necesarias para inicializar la estructura de datos que actuará como almacén de información durante el procesamiento del flujo de datos de entrada. Generalmente esto consiste en el proceso de reservar memoria, inicializar a un valor por defecto la estructura de datos, etc. Sin embargo, existen técnicas sofisticadas que requieren de una mayor carga computacional en esta fase.
- **Procesado:** Se corresponde con el procesamiento del flujo de datos de manera secuencial. La idea subyacente en esta fase es la de realizar una determinada operación sobre la estructura de datos y el elemento de entrada actual,

de manera que se lleve a cabo una actualización sobre la misma. Nótese en que la manera en que se manipula dicha estructura de datos condiciona en gran medida el conjunto de peticiones que se podrán realizar sobre ella.

- **Petición:** La fase de petición se caracteriza con respecto de la anterior por ser de carácter consultivo. Con esto nos estamos refiriendo a que dicha tarea no modifica el estado actual de la estructura de datos, sino que recoge información de la misma, que posiblemente transforme mediante alguna operación, para después obtener un valor como resultado de dicha petición.

2.4. Medidas de Análisis

Los *Algoritmos para Streaming* se caracterizan por utilizar soluciones estadísticas en alguna parte (generalmente en el procesado) de su cómputo para obtener la solución con un menor coste computacional. En este caso, el coste que se pretende minimizar es el referido al espacio necesario para almacenar la estructura de datos auxiliar. Tal y como se ha dicho anteriormente, la razón de ello es debida a que se presupone un conjunto masivo de datos en la entrada, por lo que se pretende que el orden de complejidad espacial respecto de la misma sea de carácter sublineal ($o(N)$).

El objetivo es encontrar soluciones con un índice de error acotado que permitan llegar a la solución en un orden espacial de complejidad logarítmica ($O(\log(N))$). Sin embargo, existen ocasiones en que no es posible llegar a una solución en dicho orden de complejidad, como es el caso de *Algoritmos para Streaming* aplicados a problemas de *Grafos*, en los cuales se relaja dicha restricción a un orden de complejidad *poli-logarítmico* ($O(\text{polylog}(N))$).

El orden de complejidad *poli-logarítmico* engloba el conjunto de funciones cuyo orden de complejidad presenta un crecimiento acorde a una función polinomial formada por logaritmos. Matemáticamente esto se modeliza a través de la ecuación (2.5)

$$a_k \log^k(N) + \dots + a_1 \log(N) + a_0 = O(\text{polylog}(N)) \in o(N). \quad (2.5)$$

En esta sección se muestran distintas estrategias para poder llevar a cabo la demostración de pertenencia a un determinado orden de complejidad de un *Algoritmo para Streaming*. Debido a la elevada base estadística que requieren dichas demostraciones, a continuación se definen algunos conceptos básicos en relacionadas con estimadores estadísticos, para después realizar una breve demostración acerca de distintas cotas de concentración de valores en una distribución en las subsecciones 2.4.2, 2.4.3 y 2.4.4. Las definiciones que se exponen a continuación han sido extraídas de los apuntes del curso sobre *Randomized Algorithms* [Asp16] impartido por Aspnes en la *Universidad de Yale* así como las de la asignatura de *Estadística* [SJNSB16] impartida en el Grado de Ingeniería Informática de la *Universidad de Valladolid*.

2.4.1. Conceptos básicos de Estadística

Denotaremos como x_1 a una observación cualquiera contenida en el espacio de todas las posibles. Al conjunto de todas las observaciones posibles lo denotaremos como Ω y lo denominaremos espacio muestral, por lo tanto, $x \in \Omega$. Este concepto se puede entender de manera más sencilla mediante el siguiente ejemplo. Supongamos el lanzamiento de una moneda, que como resultado puede tomar los valores cara o cruz. Definiremos entonces $x_1 = \text{cara}$ y $x_2 = \text{cruz}$ como los sucesos posibles de lanzar una moneda. Por tanto el espacio Ω se define como $\Omega = \{x_1, x_2\} = \{\text{cara}, \text{cruz}\}$.

El siguiente paso es definir el concepto de **Variable Aleatoria**, que representa una función que mapea la realización de un determinado suceso sobre el espacio Ω . Dicha función se denota con letras mayúsculas y puesto que sus parámetros de entrada son desconocidos, estos se ignoran en la notación. Por tanto denotaremos las variables aleatorias como $(\mathbf{E}, \mathbf{X}, \mathbf{Y}, \text{etc.})$. Para la variable aleatoria \mathbf{X} , sean $x_1, x_2, \dots, x_i, \dots$ cada una de las observaciones

posibles. Siguiendo el ejemplo anterior, se puede modelizar el lanzamiento de una moneda como X . Nótese por tanto, que una variable aleatoria puede definirse de manera textual como la modelización del resultado de un suceso a priori desconocido.

Definiremos probabilidad como la medida de certidumbre asociada a un suceso o evento futuro, expresada como un valor contenido en el intervalo $[0, 1]$, tomando el valor 0 un suceso imposible y 1 un suceso seguro. La notación seguida para representar esto será $Pr[\mathbf{X} = x_i]$. Suponiendo la equiprobabilidad en el ejemplo de la moneda, podemos definir sus valores de probabilidad como $Pr[\mathbf{X} = \text{cara}] = \frac{1}{2}$ y $Pr[\mathbf{X} = \text{cruz}] = \frac{1}{2}$

Una vez descritos estos conceptos simples, a continuación hablaremos sobre distintos conceptos estadísticos utilizados en el análisis de algoritmos probabilísticos tales como *Esperanza*, *Varianza*, *Variables Independientes* y *Probabilidad Condicionada*.

Denominaremos **Esperanza Matemática** al valor medio o más probable que se espera que tome una determinada variable aleatoria. La modelización matemática de dicho concepto se muestra en la ecuación (2.6). Además, la esperanza matemática es de carácter lineal, por lo que se cumplen las ecuaciones (2.7) y (2.8)

$$\mathbb{E}[\mathbf{X}] = \sum_{i=1}^{\infty} x_i \cdot Pr[\mathbf{X} = x_i] \quad (2.6)$$

$$\mathbb{E}[c\mathbf{X}] = c\mathbb{E}[\mathbf{X}] \quad (2.7)$$

$$\mathbb{E}[\mathbf{X} + \mathbf{Y}] = \mathbb{E}[\mathbf{X}] + \mathbb{E}[\mathbf{Y}] \quad (2.8)$$

La **Varianza** se define como una medida de dispersión de una variable aleatoria. Dicho estimador representa el error cuadrático respecto de la esperanza. Su modelización matemática se muestra en la ecuación (2.9). Aplicando propiedades algebraicas se puede demostrar la veracidad de las propiedades descritas en las ecuaciones (2.10) y (2.11).

$$Var[\mathbf{X}] = \mathbb{E}[(\mathbf{X} - \mathbb{E}[\mathbf{X}])^2] \quad (2.9)$$

$$Var[\mathbf{X}] = \mathbb{E}[\mathbf{X}^2] - \mathbb{E}^2[\mathbf{X}] \quad (2.10)$$

$$Var[c\mathbf{X}] = c^2 Var[\mathbf{X}] \quad (2.11)$$

A continuación se describe el concepto de **Independencia** entre dos variables aleatorias \mathbf{X}, \mathbf{Y} . Se dice que dos variables son independientes cuando los sucesos de cada una de ellas no están condicionados por los de otras. Esto puede verse a como el cumplimiento de la igualdad de la ecuación (2.12).

$$Pr[\mathbf{X} = x \cap \mathbf{Y} = y] = Pr[\mathbf{X} = x] \cdot Pr[\mathbf{Y} = y] \quad (2.12)$$

Cuando nos referimos al concepto de independencia referido a un conjunto n variables aleatorias $\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_n$ lo denominaremos **Independencia Mutua**, que impone la restricción descrita en la ecuación (2.13).

$$Pr\left[\bigcap_{i=1}^n \mathbf{X}_i = x_i\right] = \prod_{i=1}^n Pr[\mathbf{X}_i = x_i] \quad (2.13)$$

También es de especial interés en el campo de los algoritmos probabilísticos el caso de la **k-independencia** sobre un conjunto de n variables aleatorias $\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_n$. Dicha idea se puede resumir como la independencia de todas las variables en grupos de k variables. Este concepto tiene mucha importancia en el ámbito de los *Sketches*, tal y como se verá en la sección ?? . El caso más simple es para $k = 2$, el cual se denomina **independencia pareada**, cuya modelización matemática se muestra en la ecuación (2.14).

$$\forall i, \forall j \ Pr[\mathbf{X}_i = x_i \cap \mathbf{X}_j = x_j] = Pr[\mathbf{X}_i = x_i] \cdot Pr[\mathbf{X}_j = x_j] \quad (2.14)$$

Desde el punto de vista de conjuntos de n variables aleatorias $\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_n$, existen distintas propiedades de linealidad que se cumplen entre ellas a nivel del cálculo de la *Esperanza* y la *Varianza*. En el caso de la *Esperanza*, la linealidad respecto de la suma (ecuación (2.15)) se cumple para variables dependientes e independientes. Sin embargo, en el caso de la *Varianza*, la linealidad respecto de la suma (ecuación (2.16)) se cumple tan solo para variables **independientes pareadas**.

$$\mathbb{E}\left[\sum_{i=1}^n \mathbf{X}_i\right] = \sum_{i=1}^n \mathbb{E}[\mathbf{X}_i] \quad (2.15)$$

$$Var\left[\sum_{i=1}^n \mathbf{X}_i\right] = \sum_{i=1}^n Var[\mathbf{X}_i] \quad (2.16)$$

La **Probabilidad Condicionada** entre dos variables aleatorias \mathbf{E}_1 y \mathbf{E}_2 se puede definir como la medida de verosimilitud de la ocurrencia del suceso \mathbf{E}_1 sabiendo que ya ha ocurrido \mathbf{E}_2 . Esto se puede modelizar matemáticamente tal y como se muestra en la ecuación (2.17).

$$Pr[\mathbf{E}_1|\mathbf{E}_2] = \frac{Pr[\mathbf{E}_1 \cap \mathbf{E}_2]}{Pr[\mathbf{E}_2]} \quad (2.17)$$

En el caso de la **Probabilidad Condicionada** sobre variables independientes, surge la propiedad descrita en la ecuación (2.18). Es fácil entender la razón, que se apoya en la idea de que si dos variables aleatorias no guardan relación, entonces la ocurrencia de una de ellas, no condicionará el resultado de la otra.

$$Pr[\mathbf{X}_1 = x_1 | \mathbf{X}_2 = x_2] = \frac{Pr[\mathbf{X}_1 = x_1 \cap \mathbf{X}_2 = x_2]}{Pr[\mathbf{X}_2 = x_2]} = \frac{Pr[\mathbf{X}_1 = x_1] \cdot Pr[\mathbf{X}_2 = x_2]}{Pr[\mathbf{X}_2 = x_2]} = Pr[\mathbf{X}_1 = x_1] \quad (2.18)$$

Una vez descritos los conceptos estadísticos básicos para el análisis de algoritmos probabilísticos, lo siguiente es realizar una exposición acerca de las distintas cotas de concentración de valores, lo cual permite obtener resultados aproximados acerca de los resultados esperados por dichos algoritmos, así como sus niveles de complejidad. Primero se describirá *Desigualdad de Boole*, para después tratar las desigualdades de *Markov*(2.4.2), *Chebyshev*(2.4.3) y *Chernoff*(2.4.4)

La **Desigualdad de Boole** consiste en una propiedad básica que indica que la probabilidad de que se cumpla la ocurrencia de un suceso es menor o igual que ocurrencia de la suma de todas ellas. Esto se modeliza matemáticamente en la ecuación (2.19).

$$Pr\left[\bigcup_{i=1}^n \mathbf{E}_i\right] \leq \sum_{i=1}^n Pr[\mathbf{E}_i] \quad (2.19)$$

2.4.2. Desigualdad de Markov

La *Desigualdad de Markov* es la técnica base que utilizan otras desigualdades más sofisticadas para tratar de acotar la *Esperanza* de una determinada *Variable Aleatoria*. Proporciona una cota superior de probabilidad respecto

de la *Esperanza* tal y como se muestra en la ecuación (2.20). Tal y como se puede intuir, dicha cota es muy poco ajustada, sin embargo, presenta una propiedad muy interesante como estructura base. Sea $f : \mathbf{X} \rightarrow \mathbb{R}^+$ una función positiva, entonces también se cumple la desigualdad de la ecuación (2.21). El punto interesante surge cuando se escoge la función f de tal manera que sea estrictamente creciente, entonces se cumple la propiedad descrita ecuación (2.22), a través de la cual podemos obtener cotas mucho más ajustadas. Dicha idea se apoya en la *Desigualdad de Jensen*.

$$\forall \lambda \geq 0, \Pr[\mathbf{X} \geq \lambda] \leq \frac{\mathbb{E}[\mathbf{X}]}{\lambda} \quad (2.20)$$

$$\forall \lambda \geq 0, \Pr[f(\mathbf{X}) \geq f(\lambda)] \leq \frac{\mathbb{E}[f(\mathbf{X})]}{f(\lambda)} \quad (2.21)$$

$$\forall \lambda \geq 0, \Pr[\mathbf{X} \geq \lambda] = \Pr[f(\mathbf{X}) \geq f(\lambda)] \leq \frac{\mathbb{E}[f(\mathbf{X})]}{f(\lambda)} \quad (2.22)$$

2.4.3. Desigualdad de Chebyshev

La *Desigualdad de Chebyshev* utiliza la técnica descrita en la subsección anterior apoyandose en al idea de la función f para obtener una cota de concentración mucho más ajustada basandose en la *Varianza*. Dicha propiedad se muestra en la ecuación (2.23). En este caso se utiliza $f(\mathbf{X}) = \mathbf{X}^2$, que es estrictamente creciente en el dominio de aplicación de \mathbf{X} . Además, se selecciona como variable aleatoria $|\mathbf{X} - \mathbb{E}[\mathbf{X}]|$, es decir, el error absoluto de una \mathbf{X} respecto de su valor esperado. La demostración de esta idea se muestra en la ecuación (2.24).

$$\forall \lambda \geq 0, \Pr[|\mathbf{X} - \mathbb{E}[\mathbf{X}]| \geq \lambda] \leq \frac{Var[\mathbf{X}]}{\lambda^2} \quad (2.23)$$

$$\forall \lambda \geq 0, \Pr[|\mathbf{X} - \mathbb{E}[\mathbf{X}]| \geq \lambda] = \Pr[(\mathbf{X} - \mathbb{E}[\mathbf{X}])^2 \geq \lambda^2] \leq \frac{\mathbb{E}[(\mathbf{X} - \mathbb{E}[\mathbf{X}])^2]}{\lambda^2} = \frac{Var[\mathbf{X}]}{\lambda^2} \quad (2.24)$$

2.4.4. Desigualdad de Chernoff

En este apartado se realiza una descripción acerca de la *Desigualdad de Chernoff*. Dicha descripción ha sido extraída de los apuntes de la asignatura de Algoritmos Probabilísticos (*Randomized Algorithms*) [Cha04] impartida por Shuchi Chawla en la *Carnegie Mellon University* de Pennsylvania.

La *Desigualdad de Chernoff* proporciona cotas mucho más ajustadas que por contra, exigen unas presunciones más restrictivas para poder ser utilizada. La variable aleatoria en este caso debe ser de la forma $\mathbf{S} = \sum_{i=1}^n \mathbf{X}_i$ donde cada \mathbf{X}_i es una variable aleatoria uniformemente distribuida e independiente del resto. También describiremos la esperanza de cada una de las variables \mathbf{X}_i como $\mathbb{E}[\mathbf{X}_i] = p_i$.

Denotaremos como μ a la esperanza de \mathbf{S} , tal y como se describe en la ecuación (2.25). También se define la función f como $f(\mathbf{S}) = e^{t\mathbf{S}}$.

$$\mu = \mathbb{E}\left[\sum_{i=1}^n \mathbf{X}_i\right] = \sum_{i=1}^n \mathbb{E}[\mathbf{X}_i] = \sum_{i=1}^n p_i \quad (2.25)$$

El siguiente paso es utilizar la ecuación (2.22) de la *Desigualdad de Markov* con la función f , que en este caso es posible puesto que es estrictamente creciente. En este caso en lugar de utilizar λ como constante, se prefiere $\delta \in [0, 1]$, que se relaciona con la anterior de la siguiente manera: $\lambda = (1 + \delta)$. Entonces la ecuación (2.26) muestra el paso inicial para llegar a la *Desigualdad de Chernoff*.

$$Pr[\mathbf{X} > (1 + \delta)\mu] = Pr[e^{t\mathbf{X}} > e^{(1+\delta)t\mu}] \leq \frac{\mathbb{E}[e^{t\mathbf{X}}]}{e^{(1+\delta)t\mu}} \quad (2.26)$$

Aplicando operaciones aritméticas y otras propiedades estadísticas, se puede demostrar la veracidad de las ecuaciones (2.27) y (2.28), que proporcionan cotas mucho muy ajustadas de concentración de la distribución de una *variable aleatoria* formada por la suma de n variables aleatorias independientes uniformemente distribuidas.

$$\forall \delta \geq 0, Pr[\mathbf{X} \geq (1 + \delta)\mu] \leq e^{\frac{-\lambda^2 \mu}{2 + \lambda}} \quad (2.27)$$

$$\forall \delta \geq 0, Pr[\mathbf{X} \leq (1 - \delta)\mu] \leq e^{\frac{-\lambda^2 \mu}{2 + \lambda}} \quad (2.28)$$

[TODO Explicar que se van a empezar a describir algoritmos para streaming]

2.5. Algoritmo de Morris

[TODO]

2.6. Algoritmo de Flajolet-Martin

[TODO]

2.7. Aproximación a los Momentos de Frecuencia

[TODO]

Capítulo 3

Estructuras de Datos de Resumen

3.1. Introducción

[TODO]

3.2. Tipos de Estructuras de Datos de Resumen

[TODO]

3.3. Sketching

[TODO]

3.4. Count-Min Sketch

[TODO]

3.5. Count Sketch

[TODO]

3.6. AMS Sketch

[TODO]

3.7. HyperLogLog

[TODO]

Capítulo 4

Algoritmos para Grafos

4.1. Introducción

[TODO]

4.2. Modelo en Semi-Streaming

[TODO]

Capítulo 5

Reducción de la Dimensionalidad

5.1. Introducción

[TODO]

5.2. Teorema de Johnson-Lindenstrauss

[TODO]

5.3. Búsqueda de Vecinos más Cercanos

[TODO]

Capítulo 6

Técnicas de Minería de Datos

6.1. Introducción

[TODO]

6.2. Aprendizaje Supervisado y No Supervisado

[TODO]

6.3. Árboles de Decisión

[TODO]

6.4. Regresión Lineal

[TODO]

6.5. Redes Neuronales

[TODO]

6.6. Manifold Learning

[TODO]

Capítulo 7

Paralelización a Gran Escala

7.1. Introducción

[TODO]

7.2. Sistemas de Ficheros Distribuidos

[TODO]

7.3. Modelo de acceso a Memoria

[TODO]

7.4. Complejidad de la Computación

[TODO]

7.5. MapReduce

[TODO]

Apéndice A

Metodología de Trabajo

Apéndice B

¿Cómo ha sido generado este documento?

Bibliografía

- [Asp16] James Aspnes. Notes on randomized algorithms cpsc 469/569. Fall 2016.
- [Bab79] László Babai. Monte-carlo algorithms in graph isomorphism testing, 1979.
- [Cha04] Shuchi Chawla. Notes on randomized algorithms: Chernoff bounds. October 2004.
- [Ins17] Institute for Telecommunication Sciences. Definitions: Data stream, March 2017. https://www.its.bldrdoc.gov/fs-1037/dir-010/_1451.htm.
- [Kar92] Richard M. Karp. On-line algorithms versus off-line algorithms: How much is it worth to know the future? pages 416–429, 1992.
- [Mut05] S. Muthukrishnan. Data streams: Algorithms and applications. *Found. Trends Theor. Comput. Sci.*, 1(2):117–236, August 2005.
- [SJNSB16] Luis Augusto San José Nieto and Araceli Suárez Barrio. Grado en ingeniería informática: Estadística. 2016.
- [SR94] K. G. Shin and P. Ramanathan. Real-time computing: a new discipline of computer science and engineering. *Proceedings of the IEEE*, 82(1):6–24, Jan 1994.
- [Wik17] Wikipedia. Big Data, February 2017.