



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.”[2]

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
2.1.1. Computación en Tiempo Real	7
2.1.2. Problemas Dinámicos	7
2.1.3. Algoritmos Online vs Algoritmos Offline	8
2.1.4. Algoritmos Online Probabilistas vs Deterministas	8
2.2. Modelo en Streaming	8
2.2.1. Modelo de Serie Temporal	8
2.2.2. Modelo de Caja Registradora	8
2.2.3. Modelo de Molinete	8
2.3. Estructura básica	8
3. Estructuras de Datos de Resumen	9
4. Algoritmos para Grafos	11
5. Técnicas de Minería de Datos	13
6. Paralelización a Gran Escala	15
A. Metodología de Trabajo	17
B. ¿Cómo ha sido generado este documento?	19
Bibliografía	19

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 [1] 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[1] 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 *Cóputación 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[1] 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 *periodic 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*. A continuación se muestran las características de cada subgrupo:

- **Algoritmos Online:** [TODO]
- **Algoritmos Offline:** [TODO]

[TODO hablar acerca del .análisis competitivo"]

2.1.4. Algoritmos Online Probabilistas vs Deterministas

[TODO]

2.2. Modelo en Streaming

[TODO]

2.2.1. Modelo de Serie Temporal

[TODO]

2.2.2. Modelo de Caja Registradora

[TODO]

2.2.3. Modelo de Molinete

[TODO]

2.3. Estructura básica

[TODO]

Capítulo 3

Estructuras de Datos de Resumen

[TODO]

Capítulo 4

Algoritmos para Grafos

[TODO]

Capítulo 5

Técnicas de Minería de Datos

[TODO]

Capítulo 6

Paralelización a Gran Escala

[TODO]

Apéndice A

Metodología de Trabajo

Apéndice B

¿Cómo ha sido generado este documento?

Bibliografía

- [1] SHIN, K. G., AND RAMANATHAN, P. Real-time computing: a new discipline of computer science and engineering. *Proceedings of the IEEE* 82, 1 (Jan 1994), 6–24.
- [2] WIKIPEDIA. Big Data.