

PERFILADO DE NODOS
SOLUCIÓN PARA GARLANET



AUTOR: HAIZEA GÓMEZ RUIZ
DIRECTOR: JOAN MANUEL MARQUÈS PUIG
DEPARTAMENTO

CONSULTOR: ALEXANDRE VIEJO GALICIA

PERFILADO DE NODOS	1
SOLUCIÓN PARA GARLANET	1
ABSTRACT	3
INTRODUCCIÓN	3
MOTIVACIONES	4
COMPUTACIÓN APORTADA VOLUNTARIAMENTE	4
ESTADO DEL ARTE	5
SISTEMAS DE PERFILADO O CLUSTERS	5
ALGORITMOS DE PERFILADO	8
CLUSTER PROBABILISTICO	15
CALIDAD DEL CLUSTER	18
NUESTRO ANALISIS	18
METODOLOGÍA	19
ALGORITMOS SELECCIONADOS	20
RESULTADOS CON PROCEDIMIENTO PROPIO	23
ANÁLISIS CON MÉTODO DE PERFILADO K-MEANS	26
MODIFICACIÓN DEL PLANTEAMIENTO DE K-MEANS	30
DB-SCAN	35
VERIFICACIÓN	36
VERIFICACIÓN PARA EL PROCEDIMIENTO DE PERFILADO GRUESO	36
VERIFICACIÓN PARA EL PROCEDIMIENTO DE PERFILADO CON K-MEANS	36
RESULTADOS:	37
CONCLUSIONES	38
PROXIMOS PASOS:	39
SOLUCIONES COMERCIALES QUE APOYAN LA TENDENCIA	41
AGRADECIMIENTOS	44
REFERENCIAS	44
RECURSOS:	45

ABSTRACT

El estudio surge de la necesidad de perfilar la estabilidad y uptime de los nodos en un sistema de microblogging distribuido Volunteer Computing Grid. Para simular nuestros datos se han obtenido de la aplicación Seti@home, un sistema BOINC donde publica de forma masiva la información correspondiente a estas aplicaciones.

Inicialmente se realizó una taxonomía de aplicaciones similares para identificar si ya existía una tendencia en la gestión de logs para aplicaciones de microblogging.

Se establece un método de pre-filtrado de la información (gestionando los logs de inicio y fin) y generando un vector de bits de actividad para la comparación de los nodos.

En el estudio se analizan distintas tipologías de clustering y sub-familias, para finalmente verificar con tres de ellas que cumplan los requisitos (KMeans basada en particionamiento, DB-SCAN basada en densidad y un metaheurístico propio).

Los resultados, obtenidos resultan adecuados con tiempos de computación, para K-Means siempre que se utilice la matriz de vectores de uptime como features, y utilizando la PCA (Principal Component Analysis) para reducir la dimensión del vector de bits a 2 dimensiones, se ha comprobado que en esta situación el perfilado es adecuado.

Concluyendo que tecnologías aplicadas actualmente al machine learning pueden ser válidas para el planteamiento al que nos hemos enfrentado, encontrándonos la limitación actual del modelo de almacenamiento de logs que proporcionan las aplicaciones.

PALABRAS CLAVE:

CLUSTERING, VOLUNTEER COMPUTING GRID, K-MEANS, LOGS, DB-SCAN.

INTRODUCCIÓN

Actualmente, la minería de datos se ha convertido en una necesidad para determinar los datos que manejan nuestras aplicaciones.

En términos generales, podemos decir que la minería de datos es el proceso de extraer información útil, patrones y tendencias previamente desconocidas, de grandes bases de datos.

La minería de datos ha recibido un gran impulso en los últimos tiempos motivado por distintas causas:

- a) el desarrollo de algoritmos eficientes y robustos para el procesamiento de grandes volúmenes de datos,
- b) un poder computacional más barato que permite utilizar métodos analizar patrones funcionales de las aplicaciones y/o de los clientes.
- c) las ventajas comerciales y científicas que han brindado este tipo de técnicas en las más diversas áreas.

Clustering es una de las tareas más utilizadas en el proceso de minería de datos para descubrir grupos e identificar interesantes distribuciones y patrones en los datos.

En nuestro caso se pretende utilizar estas técnicas para dar una alternativa a la estabilidad de los servicios en soluciones basadas en Volunteer Computing Grid).

Garlanet, es un sistema de microblogging distribuido desarrollado por la UOC cuya pretensión es mantener la independencia de los proveedores de servicios.

Se basa en que los nodos son aportados de forma voluntaria.

El principal propósito de esta infraestructura es no tener las limitaciones de los sistemas de microblogging actuales, que son potencialmente vulnerables a ser coaccionados a:

- nivel político por entidades, gobiernos, empresas, limitando la gestión de sus contenidos,
- nivel de modelos de servicios, costes,..

Actualmente el sistema asocia el mismo grado de fiabilidad a cada uno de los nodos donde corre el servicio. Provocando en muchos casos la indisponibilidad de los mensajes que se encuentran almacenados por estar albergados en servidores menos estables.

La estabilidad de otros sistemas de microblogging distribuidos se basan en la replicación de la información en cada uno de los nodos, o la centralización de la misma en los sistemas del proveedor de servicios asociados a la aplicación de microblogging. Estas soluciones que van en contra de la filosofía de la solución.

Por lo que se pretende definir un modelado de datos que sea capaz de dar solución a la limitación actual, permitiendo analizar la aleatoriedad de los nodos y generar perfiles de candidatos aptos basado en un uptime adecuado para definirlo como candidato para almacenar los mensajes.

MOTIVACIONES

El análisis aquí planteado puede ser extrapolado a cualquier otro problema donde se pretenda calcular el perfilado de actores basado en el tiempo actividad como pudiera ser:

- El tiempo de permanencia de una identidad en una aplicación,

COMPUTACIÓN APORTADA VOLUNTARIAMENTE

La computación aportada voluntariamente (Volunteer Computing Grid o VCG), fue una iniciativa extendida por Seti@home, que junto con la universidad de Berkley desarrolló la plataforma de BOINC, para unificar y concentrar los ciclos de CPU, almacenamiento (distribuido) que aportaba la gente de forma voluntaria. Donando parte de la potencia del ordenador o dispositivo para contribuir a investigaciones científicas que requieren de un super-computador o un grid para procesar la información. Esto se basa en instalarse un programa o agente que realiza esta función.

El como y cuando se donan los recursos, depende del usuario y hacen que los nodos no sean igualmente fiables y sigan patrones distintos.

ESTADO DEL ARTE

En un primer momento, se realizó un análisis de aplicaciones similares para identificar como resolvían este problema, viendo que:

- bien utilizaban almacenamientos centralizados.
- O bien resolvían el problema, a través de la replicación del almacenamiento en cada nodo.

Este análisis queda fuera del ámbito de este estudio.

Por lo cual, en este análisis nos vamos a centrar en resolver el problema, a través de perfilar los nodos identificando aquellos que se comportan de una forma más estable, limitando el número de replicas a realizar.

Dentro de las soluciones de clustering, estas se resumen en supervisado y no supervisado, ya que el particionamiento está muy relacionado con la predicción de los siguientes valores o Machine Learning.

El clustering no supervisado está orientado a la categorización y resumir los datos. Únicamente requiere instancias (elementos a analizar con sus características).

A continuación se distinguen en una breve taxonomía de algoritmos de perfilado [7] , [8] .

SISTEMAS DE PERFILADO O CLUSTERS

Definición:

El clustering o perfilado de datos ayuda a discernir la estructura y simplifica la complejidad de cantidades masivas de datos. Es una técnica común y se utiliza en diversos campos como, aprendizaje de máquina, minería de datos, reconocimiento de patrones, análisis de imágenes y bioinformática [5] , donde la distribución de la información puede ser de cualquier tamaño y forma.

La eficiencia de los algoritmos de clustering es extremadamente necesaria cuando se trabaja con enormes bases de datos y tipos de datos de grandes dimensiones.

En [8] se realizó un análisis de los procedimientos de perfilado, que vamos a seguir.

- **Elección de las variables.**

- *Variables cualitativas:*

- Ordinales (ejemplo nivel de estudios).
- Nominales: (ejemplo nacionalidad).

- *Variables cuantitativas:*

- Variables discretas: (ejemplo nº de hijos)
- Variables continuas: (ejemplo peso).

- **Elección de la medida de asociación.**

Aquí hemos de distinguir si lo que se pretende agrupar elementos dentro de un cluster, para lo cual se ha de realizar un análisis de los elementos. O si se pretende agrupar “variables” más parecidas, para lo cual se ha de hacer un análisis del cluster de las variables, para ello basta en considerar la matriz de datos inicial X transpuesta (variando filas por columnas).

Las medidas de asociación para elementos o distancias¹:

- Distancia de Minkowski (invariante ante translaciones)

$$dq(x_i, x_j) = \|x_i - x_j\|_q = \left(\sum_{l=1}^p |x_{il} - x_{jl}|^q \right)^{1/q}; q \geq 1$$

- Distancia d1 o de ciudad (City Block) (con q=1):

$$dq(x_i, x_j) = \|x_i - x_j\| = \left(\sum_{l=1}^p |x_{il} - x_{jl}| \right)$$

- Distancia euclídea (con q=2)

$$d_2(x_i - x_j) = \|x_i - x_j\| = \sqrt{\sum_{n=1}^p (x_{in} - x_{jn})^2}$$

- Distancia de Tchebychev o del máximo (con q=infinito):

$$d_{\infty}(x_i, x_j) = \max \{l = 1, \dots, p\} |x_{il} - x_{jl}|$$

- Distancia de Mahalanobis

$$D_s(x_i, x_j) = \sqrt{(x_i - x_j)' S^{-1} (x_i - x_j)}$$

$\forall x_i \in C_i$

- **Elección de la técnica de clusters.**

Definición de particionado (jerárquico (aglomerativos, disociativos), por partición, por búsqueda de densidad), que explicaremos detenidamente en un apartado.

Existen también los métodos reductivos como el análisis factorial de tipo Q, o los métodos directos como el Block Cluster.

- **Validación de los resultados.**

Definiremos un apartado a este punto.

En nuestro caso de análisis, tenemos un conjunto de variables cuantitativas continuas, donde nuestro perfilado pretende agrupar elementos (nodos de Seti@home) en clusters.

¹ Estas distancias serán referenciadas en distintas partes del texto.

La elección del método de particionado y el método del cálculo de las distancias será definido una vez que se analicen los algoritmos de perfilado. De cara a la selección de la metodología hemos de tener en cuenta el impacto estimado de tiempo de computación en base a su eficiencia (**Error! Reference source not found.**), y el tamaño de la muestra que vamos a analizar (véase **Error! Reference source not found.**).

Tabla 1 Impacto de tiempos en base a la eficiencia

n	10	100	1.000	10.000	100.000
O(n)	10 ms	0.1s	1s	10 s	100 s
O (n*log n)	33 ms	0.7 s	10 s	2 min	28 min
O (n²)	100 ms	10 s	17 min	28 h	115 d
O (n³)	1sg	17 min	12 d	31 años	32 milenios

ALGORITMOS DE PERFILADO

Se ha analizado las siguientes familias y sub-familias de cara a identificar sus limitaciones y asociar el mecanismo adecuado acorde a nuestro volumen de información y el coste computacional:

CLUSTERS POR PARTICIONES

Se denominan algoritmos basados en disimilaridad.

Características: Se debe fijar una “k” o número de clusters.

Algoritmos más representativos:

- k-Means (descrita en el procedimiento de análisis **K-MEANS**)
- k-Modes (de 1998)
- k-Medoids
 - Es una variante utiliza medianas en el método cálculo de la distancia de los centroides para solventar el problema de outliers.
 - PAM [Partitioning Around Medoids] (de 1987)
 - CLARA [Clustering for LARge Applications] (de1990)
 - CLARANS [CLARa based on RANdomized Search] (de1994)
- BCFR [Bradley, Fayyad & Reina]
- K-prototipos: Es una alternativa en donde tenemos datos mixtos, donde por una lado se utilizará la distancia euclídea para los datos numérico y el segundo (no numérico) se usa la moda. A este segundo término se asigna un peso para poder ponderarlo y disponer de un resultado numérico.

$$d2(X, Y) = \sum_1^m \left(\sum_j^p (xi - yi)^2 + \gamma \sum_{j=i+1}^m \delta(xj, yj) \right)$$

Amplía la posibilidad de usar K-Modas para datos categóricos, y utilizar una solución similar a K-Means.

K-MEANS

Posiblemente sea el algoritmo de particiones más conocido y testado [4] debido a que comparativamente a aportado los mejores rendimientos basados en pruebas y por tanto seleccionado para nuestro análisis.

Es un método que es sensible a la presencia de outliers, por lo que se precisa un pre-filltrado de estos elementos [D.1].

Este método se basa en hacer perfilados conociendo el número de clusters (k).

Cada cluster tiene asociado un centroide (centro geométrico del cluster). Los puntos se asignan al cluster cuyo centroide esté más cerca (utilizando cualquier métrica de distancia).

Iterativamente, se van actualizando los centroides en función de las asignaciones de puntos a clusters, hasta que los centroides dejen de cambiar.

Limitaciones:

Existen limitaciones adicionales para la implementación de esta solución:

- Cluster no convexos, de distinto tamaño o de distinta densidad.
- Selección adecuada de los centroides.
- Selección del algoritmo para el cálculo de centroides.
- Atributos a comparar no son numéricos.

Soluciones:

Solución para clusters no convexos, de distinto tamaño o de distinta densidad:

- Se implementa K-Means con valores de “k” elevados y se revisan los resultados obtenidos.
- Se implementan soluciones alternativas como “Spectral clustering”

Soluciones para la selección de centroides:

1 – Partir de centroides escogidos al azar, pero realizar las suficientes iteraciones de medición de la distancia entre nodos, seleccionando los nuevos centroides y comparar los resultados obtenidos (GRASP) . Metodología implementada en nuestro análisis.

2 – Escoger centroides de partida potenciales (pre-calculados), para ello existen distintos métodos (k-means++) implementado en [D.2] .

3- Seleccionar un método jerárquico para determinar el valor de “k” y posteriormente implementar un sistema basado en particiones.

Solución para comparar atributos que no son numéricos:

Nota: Aunque K-Means es un algoritmo válido para strings (usando la distancia de Levenshtein) que es similar a la distancia de Hamming (binaria), aquí nos referimos a datos no estructurados.

Nota: La similaridad entre dos cadenas de texto A y B se basa en el conjunto mínimo de operaciones de edición necesarias para transformar “A” en “B”, o viceversa. Hay tres operaciones de edición, las cuales son destrucción, inserción y substitución. Entre más cerca de cero es la distancia de Levenshtein más parecidas son las hileras.

Nota: La forma de recalculer los centroides:

$$c_j = \frac{1}{|C_j|} * \sum_{z \in C_j} z$$

donde “c_j” es el centroide y “z” es un elemento que pertenece al cluster “C_j”.

Se pretende minimizar las distancias inter-cluster o separación entre clusters e intra-cluster o cohesión de los nodos.

Algunas definiciones con respecto a las distancias:

La distancia promedio inter-clústeres es la distancia promedio entre los centroides de los clusters.

- La distancia mínima inter-clusters (centroides) es la distancia mínima entre dos centroides de dos clusters.
- La distancia mínima entre clusters (nodos) es la distancia mínima entre un nodo y el centroide de un clúster al que no pertenece.
- La distancia promedio dentro del clúster (intra-cluster) es la distancia promedio entre un nodo y el centroide del clúster al que pertenece.
- La distancia máxima intra-clúster es la distancia máxima entre un nodo y el centroide del clúster al que pertenece.

Complejidad Tiempo: $O(n * k * I * d)$

Donde n = número de puntos, k = número de clusters, I = número iteraciones, d = número de atributos

Complejidad en el espacio: $O((N + K) * D)$

Comparativa de eficiencia entre métodos:

PAM, limitado a su complejidad computacional (escalabilidad), es válido para conjunto de datos pequeños.

Comparativa de eficiencia:

PAM $O(k * (n - k)^2)$ vs K-Means $O(n * k)$

MÉTODOS BASADOS EN EL ANALISIS DE DENSIDAD

Este tipo de análisis es efectivo cuando su forma es irregular, están entrelazados o existe ruido y outliers de datos.

Pretendemos verificar que nuestra suposición de vectores y características es compatible con el algoritmo de densidad Db-scan.

Características:

- Se identifican a los clusters de forma arbitraria.

- Robustos ante la presencia de ruido (puntos dispersos que no forman parte de un cluster ni están en su frontera) .
- Escalables: Un único recorrido del conjunto de datos.

Algoritmos más comunes:

- DBSCAN: Density Based Spatial Cluster Applications with Noise (KDD 1996)
- OPTICS: Ordering Points To Identify The Cluster Structure (Ankerst et al SIGMOD 1999)
- DENCURE: Density based Clustering (Hinnedburg & Keim, KDD 1998).

DBSCAN:

Utilizaremos este algoritmo como comprobación [D.4] .

Es el más extendido de los clusters de densidad por su eficiencia y facilidad ya que requiere pocos parámetros de entrada, para establecer una partición:

Elimina el ruido y agrupa los datos restantes:

Se han de determinar los parámetros :

- Epsilon (radio del vecindario o del cluster).
- minPts (mínimo número de puntos del vecindario o cluster).

Eficiencia tiempo: $O(n * \text{tiempo necesario para encontrar vecinos en la distancia Epsilon})$

En el peor caso: $O(n^2)$

En el mejor caso: $O(n * \log n)$.

Eficiencia en espacio: $O(n^2)$

Limitaciones:

No es capaz de detectar clusters anidados.

Solución:

Implementar la alternativa de **enDBSCAN** (modificación de DBSCAN capaz de detectar cluster anidados en espacios de densidad variable).

Se agrupan en base a su densidad y el radio prefijado.

CLUSTERING JERÁRQUICOS

Características: No se fija la k .

Algoritmos más representativos:

- BIRCH: Balanced Iterative Reducing and Clustering using Hierarchies (Zhang, Ramakrishnan & Livny, SIGMOD'1996)
- CURE: Clustering Using Representatives (Guha, Rastogi & Shim, SIGMOD'1998)
- ROCK: Robust Clustering using links (Guha, Rastogi & Shim, ICDE'1999)
- CHAMELEON: Hierarchical Clustering Using Dynamic Modeling (Karypis, Han & Kumar, 1999)

Son métodos de agrupamiento que no requiere fijar la " k ", no necesitan eliminar outliers y no están limitados a la comparación de clusters de distintas densidades. Pero por contrapartida tienen una complejidad de orden superiores.

Una vez que se conocen las distancias existentes entre cada dos elementos se observa cuáles son los individuos más próximos en cuanto a esta distancia o similitud. Estos formarán un grupo que no vuelve a separarse durante el proceso. Se mide la distancia o similitud entre todos los elementos de nuevo (tomando el grupo como un elemento) de la siguiente forma:

- Cuando se mide la distancia entre el grupo formado y un elemento, se toma la distancia mínima de los elementos del grupo al nuevo elemento,
- Cuando se mide la similitud o similaridad entre el grupo formado y un elemento, se toma la máxima de los elementos del grupo al nuevo elemento.

Distancia	A	B	C	D
A	0			
B	9	0		
C	4	5		
D	7	3	11	0

Tabla 2 Ejemplo de matriz de similitudes

Distancia mínima:

Entre B y D (3) -> B y D forman un grupo, valor de B-D(7) máximo que da D.

Se miden las distancias de nuevo:

Nueva distancia mínima: C y A(4) -> A-C conforman un grupo (valor 5).

Se unen los grupos A-C y B-D formando el grupo A-B-C-D.

Aunque matemáticamente tienen una complejidad mayor y se desaconsejan para grandes volúmenes de datos se han utilizado de forma correcta en análisis como perfilado en el análisis de datos orientado a la diabetes (gran volumen de información), análisis neuronales,....

Estos métodos de forma gráfica se puede representar en forma de dendogramas o diagramas de datos basados en árboles. Agrupando los objetos por su representante y su similitud a través de la altura.

El dendograma es un gráfico usado que permite visualizar el procedimiento de agrupamiento del cluster por pasos, y que ayuda a decidir el número de grupos que podría representar la mejor estructura de datos teniendo en cuenta la forma en que se va generando el cluster y la media de similitud y las particiones que se van generando.

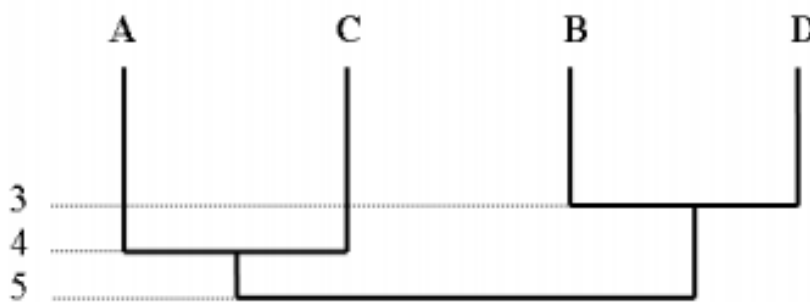


Figure 1 Dendograma

Aquí partimos que lo que nos interesa es saber el número de clusters (valor "k" de K-Means) óptimo.

Siempre que utilicemos una muestra baja, ya que el coste computacional es alto.

Este método nos identifica los outliers de forma que no se requiere pre-procesar los datos.

Aquí para agrupar los datos podemos seguir dos tendencias de clasificación:

- Técnicas Divisivas (**DIANA**: DIvisive ANAlisys), partiendo de un conjunto de datos (todos los datos) vamos de forma recursiva generando subconjuntos de datos cohesionados llegando a un cluster o a "n" clusters.
- Técnicas Aglomerativas (**AGNES**: AGLomerative NESTing), de un subconjunto de datos perfilados "clusters" componemos el árbol de

relaciones.

Su escalabilidad es del orden:

Complejidad en el espacio: $O > (n^2)$,

Complejidad en el tiempo: $O (n^2 * \log n)$

para solventar este problema, se escogen variantes que vamos a describir:

BRICH se basa en un árbol balanceado (cuyos nodos almacenan las sumas de los Clusters Features (CF) de sus descendientes.

Es óptimo para grandes tamaños porque su escalabilidad es lineal $O (n)$.

Limitaciones:

- Sólo es válido para datos numéricos.
- Se debe fijar a priori el tamaño de clusters (para poder dividirlo en ramas) y el cluster pueden resultar no naturales

Generan una estructura de datos jerárquica CF Tree,

Los CFs tienen dos parámetros, el primero el número máximo de hijos de sus sub-nodos y el segundo el diámetro máximo de los sub-clusters. Con ambos parámetros determinados la densidad y tamaño del cluster y el peso que tiene el mismo.

Se debe determinar el número máximo del diámetro del cluster para generar una agrupación adecuada.

La variante más utilizada en las investigaciones para los métodos jerárquicos es **CURE** (Clustering Uses Representative) al que no le afectan clusters con formas convexas.

Al igual que K-Means utiliza un centroide calculado con la distancia euclídea² entre puntos, al cual le añade la información de puntos representativos, para poder ser inmune a la forma del cluster. Estos valores dispersos que serían outliers, se les acerca al centroide una distancia marcada por % definido para que puedan incluirse en el cluster en vez de prescindir de ellos.

La complejidad del algoritmo es para un espacio de orden $O (n)$, tiene un tiempo de computación del orden $O (n^2 * \log n)$.

Esto nos limita el número de elementos a procesar. Esto es debido que este tipo de clusters se implementan en memoria principal.

² Distancia euclídea: $D_{ij} = \sqrt{(\sum_{k=1}^n (x_{ki} - x_{kj})^2)}$

OTROS ALGORITMOS ANALIZADOS:

Con el fin de que no parezcan que las soluciones únicamente se basan en las anteriormente indicadas. [2]

Por mencionar algunos adicionales:

- Métodos basados en rejillas.
- Clústeres probabilísticos
- Métodos basados en Co-ocurrencia de datos categóricos.
- Clúster basados en restricciones.
- Clúster sub-espacial.

CLUSTER PROBABILISTICO

Vamos a hacer una referencia como alternativa a nuestra solución.

Existen algoritmos basados en modelos, en estos los datos son generados por modelos probabilísticos (cadenas de Markov o variantes de estas) y estimar sus parámetros. [8]

Esto le aporta una gestión de dinamismo a los análisis.

La idea de las propuestas de Cadez et al. (2003) y Pamminger and Frühwirth-Schnatter (2010), suponiendo el primero un modelo basado en cadenas de Markov para definir el comportamiento de los usuarios de la web www.msnbc.com a través de las distintas categorías de información que ofrece la página (actualidad, deportes, tiempo, viajes, página principal, ...). Para ello, gracias a las cookies de cada usuario y por medio de una transformación de los datos obtenidos a través de ellas, se ha simplificado la información pasando a obtener una secuencia de entradas a categorías por cada uno de los visitantes de la web.

Lo importante es que se supone que cada cluster va seguir un modelo de Markov distinto. Y esto aporta el planteamiento consigue indicar que la heterogeneidad existente en los distintos clusters debida a la diversidad de la duración de la estancia en cada categoría, los contenidos visitados, . . . los datos tratados en los ejemplos del artículo son modelizados como si se hubiesen generado de la siguiente manera.

- 1) A cada nodo se le asigna a un cluster con cierta probabilidad.
- 2) El comportamiento del nodo genera a partir de un modelo de Markov con parámetros específicos de ese grupo.

Así pues, debe de estar concretado el comportamiento de los nodos, para cada uno de los clusters. para determinar la proporción de nodos existentes en cada cluster y especificar los parámetros de cada modelo de Markov, se utiliza el algoritmo EM.

EM

Pertenece a una familia de modelos que se conocen como Finite Mixture Models, los cuales se pueden utilizar para segmentar conjuntos de datos.

Dentro de la clasificación, encajaría dentro de los Métodos de Particionado y Recolocación, en concreto Clustering Probabilístico. Se trata de obtener la FDP (Función de Densidad de Probabilidad) desconocida a la que pertenecen el conjunto completo de datos.

COBWEB

Este algoritmo lo mencionamos ya que la comprobación de nuestros modelos de clusters la hemos sacado de cómo COBWeb calcula el valor de Acuity, similar a como nosotros hemos reducido las características con PCA.

COBWEB [5] se basa en un cluster jerárquico y se caracteriza porque utiliza aprendizaje incremental, esto es, realiza las particiones instancia a instancia. Al algoritmo COBWEB no hay que proporcionarle el número exacto de clusters que queremos, sino que en base a los parámetros anteriormente mencionados encuentra el número óptimo

Al principio, el árbol consiste en un único nodo raíz y las instancias se van añadiendo una a una al árbol, y este se va actualizando.

La actualización consiste en encontrar el mejor sitio donde incluir la nueva instancia, operación que puede necesitar de la reestructuración de todo el árbol (incluyendo la generación de un nuevo nodo anfitrión para la instancia y/o la fusión/partición de nodos existentes) o simplemente la inclusión de la instancia en un nodo que ya existía. La clave para saber cómo y dónde se debe actualizar el árbol la proporciona una medida denominada utilidad de categoría, que mide la calidad general de una partición de instancias en un segmento.

La reestructuración que mayor utilidad de categoría proporcione es la que se adopta en ese paso.

El algoritmo es muy sensible a otros dos parámetros:

- Acuity Este parámetro es muy necesario, ya que la utilidad de categoría se basa en una estimación de la media y la desviación estándar del valor de los atributos, pero cuando se estima la desviación estándar del valor de un atributo para un nodo en particular, el resultado es cero si dicho nodo sólo contiene una instancia. Así pues, el parámetro acuity representa la medida de error de un nodo con una sola instancia, es decir, establece la varianza mínima de un atributo.
- Cut-off Este valor se utiliza para evitar el crecimiento desmesurado del número de segmentos. Indica el grado de mejoría que se debe producir en la utilidad de categoría para que la instancia sea tenida en cuenta de manera individual. En otras palabras: cuando no es suficiente el incremento de la utilidad de categoría en el momento en el que se añade un nuevo nodo, ese nodo se corta, conteniendo la instancia otro nodo ya existente.

De cara a nuestro análisis se basa en usar un conjunto de parámetros “t” que nos identifican a nuestros clusters, sin requerir utilizar todos los valores.

Y presentando una complejidad de $O(k*n*t)$ con k =número de clusters, y “ n ” número de elementos, mejorando los tiempos de ejecución de los modelos de particionamiento.

El planteamiento se basa en asociar probabilísticamente un cluster en función del %average medio del cluúster,y su desviación estándar.

Una vez determinado el cluster, su comportamiento viene definido a través de un modelo estadístico (cadena de Markov) propio del cluster al que pertenece. Con estas suposiciones acerca de la estructura, los objetivos del algoritmo son calcular el número de clusters, hallar la función de probabilidad que asigna un vector de uptime a un cluster en particular y determinar los parámetros de cada cluster.

Suponiendo que esta composición (denominada estructura de modelos mixtos) consta de k conglomerados, se puede considerar que la probabilidad de que conocidos los parámetros θ del modelo, se obtenga el elemento X viene determinada por:

$$p(X|\theta) = \sum_{i=1}^k p(c_i|\theta) p_i(X|c_i, \theta)$$

- $X = [x_1, \dots, x_m]$ representa un vector de bits. Así, X determina un elemento del conjunto de datos, con valores concretos x_i para cada una de las categorías. Siendo $i \in \{1, \dots, m\}$.
- θ hace referencia al conjunto de parámetros que describen el modelo considerado. los $c_l : l \in \{1, \dots, k\}$ se consideran indicadores de cada uno de los k clúster supuestos para el modelo.
- $p(c_l | \theta)$ es la probabilidad marginal de que un elemento pertenezca al “ i – ésimo” cluster. Como consecuencia inmediata se tiene que $\sum_{i=1}^k p(c_i|\theta) = 1$
- $p_l(X | c_i, \theta)$ es el modelo estadístico que describe la distribución para las variables dentro del “ i – ésimo” cluster; es decir la probabilidad de obtener el vector de bist X sabiendo que pertenece al cluster especificado.

CALIDAD DEL CLUSTER

A efectos de medir la calidad de los clusters encontrados por los distintos algoritmos anteriormente descritos, hemos considerado las medidas intra-cluster e inter-cluster, que nos determinan la cohesión y la separación entre cluster respectivamente.

Vamos a describir las formulas que implementaremos :

Cuantificación del error:

$$1 \text{ Error - cluster} = \frac{\sum_{i=1}^K (\sum d^{(ci-Z)})}{K}$$

Distancia Intra-Cluster:

$$2 \text{ Intra - Cluster} = \frac{\sum_{i=1}^K (\sum_{z,t \in ci} \frac{d(z,t)}{|ci|})}{K}$$

Distancia inter-Cluster:

$$3 \text{ Inter - cluster} = \frac{\sum_{i=1}^{K-1} \sum_{j=i+1}^K d(ci,cj)}{\sum_{i=1}^{K-1} i}$$

NUESTRO ANALISIS

Dado que no se dispone de logs reales para el análisis se ha propuesto el uso de los logs que aporta Seti@home que son compatibles con la solución de Garlanet.

De cara al análisis de los datos, se pretende perfilar los nodos en base a su actividad.

En nuestros casos vamos a partir de un conjunto de una cantidad de información a procesar muy elevada que nos propone su gestión basada en el análisis de fichero, o basado en queries a una base de datos relacional (mysql).

El fichero de logs análisis de partida extraído de la página FTA (Failure Trace Archive, http://fta.scem.uws.edu.au/data/seti09/seti09_tab.tgz) nos proporciona un archivo comprimido en formatos tar.gz para su análisis comprendidos entre las fechas 2007-2009.

De entre los ficheros o tablas que nos propone la FTA, se selecciona el fichero event_trace.tab cortado, ya que el original es de 18Gb de información.

Este fichero dispone de la información relativa a los estados de actividad.

Tabla 3 Ejemplo de tabla de logs event_trace.tab

event_id	component_id	node_id	platform_id	node_name	event_type	event_start_time	event_stop_time	event_end_reason
0	0	100416828	2	100416828	1	1211524407.906	1211524417.922	NULL
1	0	100416828	2	100416828	0	1211524417.922	1211524440.312	NULL
2	0	100416828	2	100416828	1	1211524440.312	1211615915.156	NULL
3	0	100416828	2	100416828	0	1211615915.156	1211616057.031	NULL

Descripción de los campos relevantes:

- Event_type: Campo que nos va a identificar si estamos activos o inactivos.
- Event_start_time: Hora de actividad en formato epoch.
- Event_stop_time: Hora de parada en formato epoch.

Dado el tamaño de nuestra muestra, se ha escogido una porción de esta información que comprende el periodo de tiempo de:

Inicio: 08:37:27 Viernes 23 de mayo de 2008.

Fin: 09:50:46, Domingo 28 de septiembre de 2008

Que corresponden a 10.000.000 de líneas de logs con un tamaño de 893MB.

De este fichero base de partida, se han extraído distintas muestras para realizar los ensayos con las siguientes características:

Tabla 4 Tamaños de las muestras disponibles

Líneas procesadas	Tamaño	Nº de nodos distintos
1.000	89 (KB)	4
10.000	889 (KB)	30
100.000	8,7 (MB)	118
1.000.000	87 (MB)	1.093
10.000.000	893 (MB)	11.236
202.546.161	18 (GB)	226.209

A través de esta información se puede deducir que el comportamiento de los nodos que continuamente se van a conectar y desconectar.

Esto se aprecia

METODOLOGÍA

En el procedimiento a realizar, pretendemos cotejar dos métodos de perfilado.

En primer lugar acondicionamos las horas de proceso acumulativas de cada nodo obteniendo la media de uso de cada hora.

A cada hora se le va a asociar un bit que identifica:

- O el average de uso de esa hora (valor para las soluciones de K-MEANS, DB-SCAN).
- O un valor binario para el algoritmo propio.

Para la selección del número de bits se ha tenido en cuenta (el coste computacional que se requeriría y la necesidad de histórico del perfilado). La elección de 24 bits correspondiente a un día o 48 bits correspondientes a 2 días, nos parecía escaso y disponer de 744 bits (correspondiente a un mes) nos pareció excesivo de cara a las continuas comparaciones que requieren los valores.

Nota: Se podría haber escogido una medida distinta como “medias horas” para tener un mejor ajuste de tiempo.

Por tanto, se ha escogido una semana 168 bits de forma que pueda cotejarse la funcionalidad y ampliar en base a los resultados un posible perfilado de un mes de histórico.

El primero de ellos, denominado perfilado propio, la queremos cotejar con el uso los métodos frecuentes de perfilado (particiones, densidad, jerárquicos), descritos anteriormente.

De cara a nuestra solución podría valer más de una método de perfilado.

En nuestro caso se ha seleccionado el método de **K-MEANS** por tener un menor coste computacional, por ser el más implantado y las comparativas de rendimiento y efectividad son superiores.

No nos vemos afectados a las limitaciones anteriormente descritas, dado que conocemos a priori el número "k" y verificamos que los clusters que obtenemos tras una primera prueba, no forman una estructura de cluster convexa y podemos eliminar los outliers previo a la tarea de procesado. Una vez definido el método de particionado, realizamos distintas pruebas para obtener la escala de datos idónea, el valor de "k" (número de clusters), y el número de iteraciones.

De igual forma que el caso de **procedimiento propio**, se pretende extraer un log y comprobar que se perfila de forma similar y con mejores rendimientos.

Para describir el ámbito de requerimientos computacionales (ya que se va a requerir varias iteraciones de comparación entre nodos), se indicará el orden de complejidad de las metodologías utilizadas) para que posteriormente pueda hacerse una similitud con otras investigaciones.

ALGORITMOS SELECCIONADOS

DEFINICION PROCEDIMIENTO PROPIO

El procedimiento propio, es el nombre que hemos asociado a nuestra solución meta heurística³.

Se basa en una simplificación de los procedimientos [3] de perfilado jerárquico, en donde no vamos a tener que mantener los datos acumulativos para gestionar los clusters y una comparativa adicional de los datos acumulativos con respecto a los nuevos puntos.

Pretendemos bajar el tiempo de computación: a una solución del orden:

$$O(n * \log n)$$

³ Meta heurística es la combinación de varios métodos heurísticos a alto nivel, para conseguir una búsqueda en el espacio de búsqueda propio, que pueda optimizar los métodos generalistas.

Estos algoritmos son aproximativos (no garantizan la obtención de la solución óptima, son relativamente sencillos, y la regla de selección se basa en el instante del proceso y de la historia hasta ese momento)

En el proceso de pre-procesado de los datos se ha establecido un vector de bits con la transformación binaria del %average de funcionamiento. Estableciendo un a 1 si el nodo ha estado activo (con datos acumulativos) más del 75% de la hora y definiendo como 0 en el caso contrario (simplificación).

Esto genera un vector de bits de 168 posiciones para cada uno de los nodos. De forma sencilla podemos establecer que nodos están continuamente activos y funcionales (con una tasa de actividad por hora superior al 75%) y cuales no.

$$actividad = \sum_1^{168} x_i$$

Entre cada una de los perfiles, nos podemos encontrar que para un mismo número de bits activos, ejemplo entre 90 a 100 horas de actividad al mes, podemos encontrarnos que esta situación se puede dar en nodos que estén:

- 3 horas por la mañana activos durante el mes.
- 3 horas nocturnas durante el mes.
- 90 horas continuas de actividad...

A través de esta solución, podemos contar las horas de actividad y con ello asociar estos nodos a diferentes perfiles de funcionamiento, esto nos lo determinará el número de horas (x) de cada array de 168 posiciones.

Se ha establecido un rango de bits que definen los perfiles en base a horas de actividad.

Tabla 5 Perfiles asociados al método propio

Valor de estabilidad	Perfil	Rango
x	Estable	$121 \leq x < \leq 168$
	Semi Estable	$81 \leq x \leq 120$
	Poco Estable	$21 \leq x \leq 80$
	Inestable	$20 \leq x$

Nota: Estos rangos de funcionamiento se han dejado parametrizados para poder variarse en valor o eliminarse perfiles.

Nota: Los valores de continuamente activo e inactivo, se tendrán en cuenta para la eliminación de outliers del proceso de K-Means.

Usaremos la distancia de Hamming para ver como un nodo se parece más o menos en base a su actividad, pudiendo agrupar estos nodos que tienen valores pequeños en la distancia de Hamming en un nuevo sub-perfil.

Se han establecido realizar el cálculo de las distancias de Hamming de forma independiente a cada perfil de perfilado, ya que mejora sustancialmente los rendimientos de cálculo.

Donde estamos comparando sólo los elementos del mismo perfil, siendo el peor de los casos, cuando todos los nodos se encuentran en el mismo perfil.

A demás nos permite comparar las horas de discrepancia concretas entre nodos semejantes mediante los valores de offset y size.

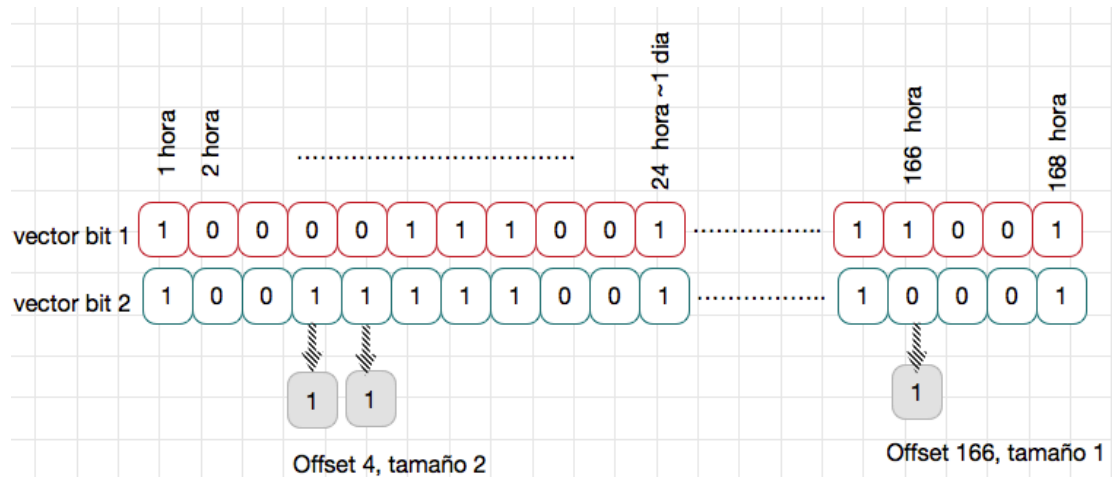


Figure 2 Método de comparación de vectores con la distancia de Hamming

Existe una alternativa que puede mejorar el perfilado, se basa en no hacer la transformación binaria, o simplificación booleana (requerida por Hamming). En este caso truncaríamos el %average a dos decimales que correspondiente al su valor real de uso.

Estos valores son traducidos a cadenas de caracteres (string casting), e implementar una comparación de caracteres. Esto lo implementaríamos con la distancia de *Levenshtein*.

Esta metodología nos haría definir un pre-análisis para determinar los outliers, pero en contrapartida nos mostraría aproximaciones más finas.

100203236	0			
100422926	76	0		
100413314	81	71	0	
100309685	79	73	100	0

Apreciamos que los valores de similitud en los nodos poco estables o semi estables son muy elevados, por lo que hay poca cohesión en el grupo. No así en los nodos estables.

Con otro tipo de agrupación (por días) o variando los umbrales de las tipologías podríamos adecuar estos valores.

A modo de ejemplo se muestra para el entorno semi-estable en las primeras 72 horas (esto saldría de forma directa variando el tamaño del array a gestionar de 168 a 72):

Tabla 9 Matriz de similitudes de nodos semi-estables en las primeras 72 horas

Nodo semi estables				
	100203236	100422926	100413314	100309685
100203236	0			
100422926	37	0		
100413314	39	23	0	
100309685	41	28	37	0

Nota: Opcionalmente: se puede realizar la “Matriz de similitudes por día”.

Vemos que aunque la actividad global de cada uno de los nodos se mantienen dentro del rango descrito (81 a 120 horas al mes), durante los primeros 3 días, las horas de actividad varían un 50% de las veces.

La gestión de (crear, almacenar, actualizar) las matrices de similitud para calcular el vecino más próximo conlleva unos costes computacionales altos cuando el número de elementos es alto (de ahí nuestra segmentación previa en perfiles), pero cuando esto ocurre se suelen emplear otros métodos de agrupamiento basado en los algoritmos de particionamiento.

Calculamos nuestros tiempos de computación:

Tabla 10 Comparaciones y tiempos de computación

Líneas de logs procesadas	Comparaciones requeridas	Tiempo de computación
1.000	21	0,08 segundos
10.000	91 / 34 optimizado	0,6 segundos
100.000	739	11 segundos

Vamos a comparar esta metodología r con los resultados que nos propone el método K-Means para verificar que se cumple y que podemos utilizar el perfilado a medida.

ANÁLISIS CON MÉTODO DE PERFILADO K-MEANS

Para este análisis como vamos a utilizar el método de particionamiento de KMeans [11], hemos filtrado y eliminados los valores acumulativos en una hora superiores (valore de $y > 0,98$) por considerarlo un nodo muy estable (siempre activo), pudiéndose haber conectado y desconectado varias veces en la misma hora, pero de forma global activo.

De igual forma se han eliminado aquellos nodos cuya actividad acumulada en una hora no ha sido superior a 0.01%, por suponerse como nodo inactivo (barras oscuras).

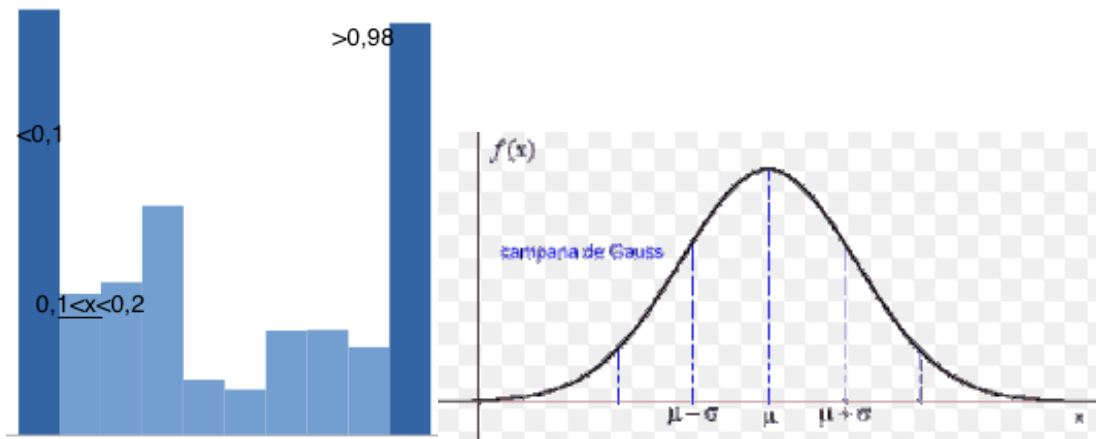


Figure 3 Densidad de %avg para 10.000 líneas de logs

Con un valor de media: 0,47

Varianza: 0,14

Desviación estándar: 0,38

No se ha realizado un ajuste o simplificación de datos a un entorno binario

Para el agrupamiento en K-Means se ha desarrollado un script en Python:

6-create-kmeans.py :

Genera los clusters, a través de un fichero de entrada con los valores de los puntos siendo:

X: hora de la medida (array de 168 posiciones).

Y: %average de actividad del nodo.

Genera los ficheros de logs de los datos (de cara a las comprobaciones).

Dibuja los gráficos de los puntos y los clusters.

Nota: Se ha utilizado Python por su capacidad de utilizar librerías estadísticas y de perfilados de ejemplo.

Realizando esto mismo, con otros centroides escogidos por valores al azar:

Valores de partida	k	C0	C1	C3	C4	Centroides finales
Numero de puntos: 63.144	2	elementos	elementos	elementos	elementos	[1,67 0.094] [1.669 0.647]

Centroides iniciales	[0.75 0.96] [0.45 0.168]		31.572	31.572	-	-	
Estructura	Lineal		y < 0.922	y > 0.922	-	-	
Otros datos			Media ⁴ : 0,667	Media: 0.875			
			Varianza: ⁵ 0,228	Varianza: 0.237			
			Desviación estándar ⁶ : 0,47756	Desviación estándar: 0,4866			
Valores de partida			elementos	elementos	elementos		
Número de puntos: 63.144		3	22.377	22.421	18.706		[1,67 0,094] [1.489 0.99] [1,669 0.647]
Centroide iniciales	[0.75 0.96] [0.45 0.168] [0.167 0.55] [0.167 0.55]						
Estructura	Lineal		y>0.3	y> 0.94	0.3<y<0.9 4	-	
			Elementos	Elementos	Elementos	Elementos	
Número de puntos: 63.144		4	14.876	14.604	15.854	18.170	[1,48 0,99] [1.66 0.64] [1.67 0.96] [1.167 0.097]
Centroides iniciales	[0.75 0.96] [0.45 0.168] [0.127 0.55] [0.167 0.88]						
Estructura	Lineal		(1,21<x<1,6 8, 0,4<y<0,6) (0.01<x<0,1 68, 0,0<y<9,16)	0.6<x, y<0,95	(0<x, y~0,34) (0,4 >x, y~0,97)	(y > 0.98)	

Tabla 11 Resultados iniciales del proceso de particionado con K-Means

⁴ Media aritmética: $\bar{x} = \frac{1}{n} * \sum_{i=1}^n x_i$

⁵ Varianza: a qué distancia los datos están agrupados alrededor de la media.

$$\sigma_n^2 = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^2$$

⁶ Desviación estándar: cuán dispersos están los números en la muestra de

datos. $s = \sqrt{\sigma_n^2} = \sqrt{\frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^2}$

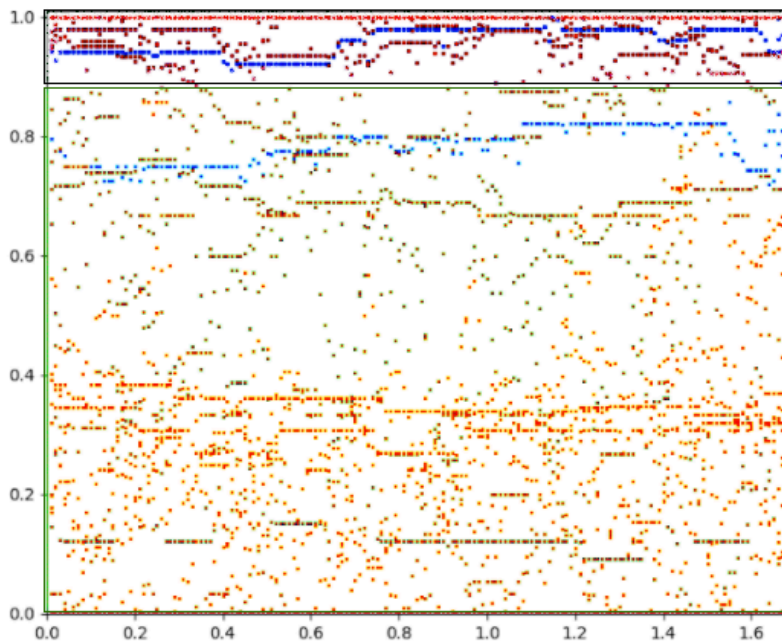


Figure 4 KMeans perfilado de logs para k=2

Nota: para k=2, el sistema se comporta como el análisis de perfilado con comparación basada en la distancia de Hamming, pudiéndose haber usando la equivalencia de 0, 1 a activo o inactivo. En este caso, si el porcentaje de actividad de un nodo en una hora determinada está cercana al valor inferior marcado por el cluster 0.97 estará en ese cluster y si su actividad media en esa hora es superior a 0.922 estará en el cluster 1. Lo normal es que un nodo tenga una comportamiento similar. En caso de no ser así, pasaría a formar parte del cluster0 o cluster 1 en función de su tanto % de uptime.

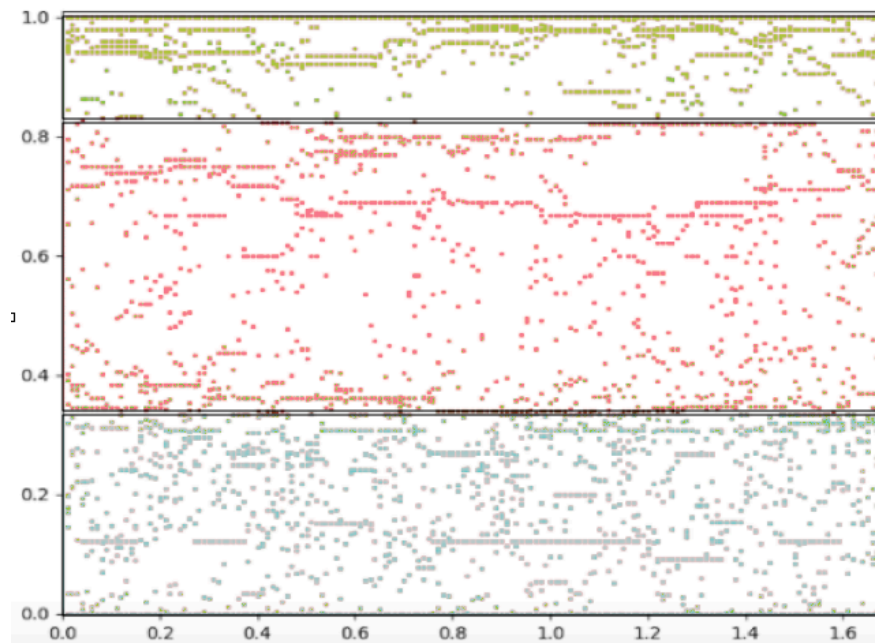


Figure 5 K-Means clústeres logs k=3, agrupado

Nota: El cluster 3 tiene menores intervalos de disponibilidad en la media (solo el 10% de los hosts tienen intervalos de disponibilidad inferiores a 10 horas en promedio), pero la desviación estándar⁷ es mucho menor que la que muestra el cluster 2. El cluster 1, tiene los mayores valores de disponibilidad y con una desviación estándar más baja.

En el proceso de verificación descrito en [13] se describe el proceso de análisis de los valores de los nodos que se encuentran en clusters distintos y que nos lleva a la deducción de dos falsos positivos que invalidan nuestro planteamiento.

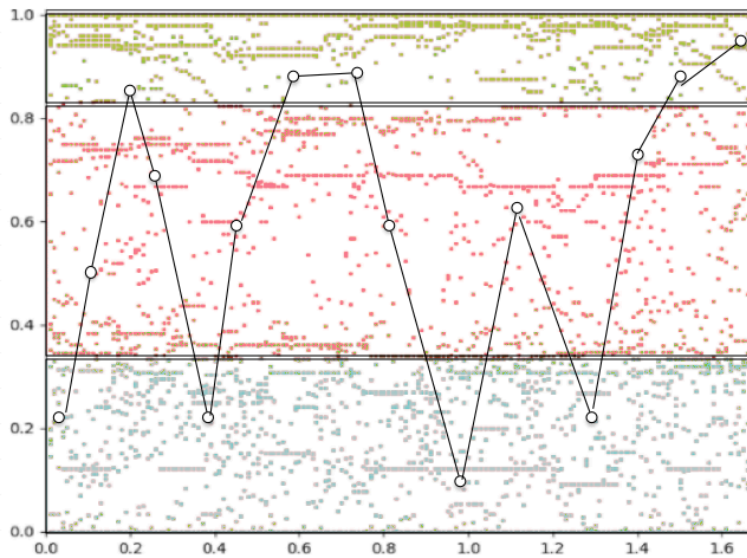


Figure 6 Potencial falso positivo¹ con un nodo con comportamiento muy aleatorio

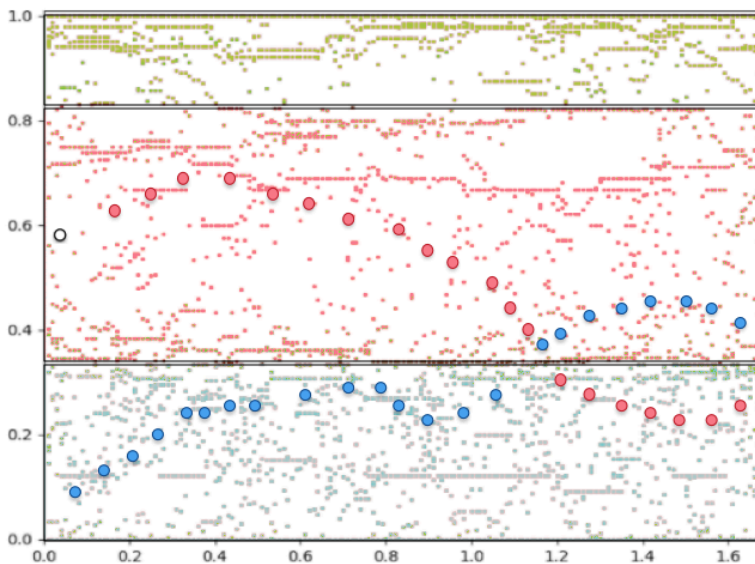


Figure 7 Posible falso positivo de nodos con tendencias contrarias y puntos cercanos al límite del cluster

⁷ Desviación estándar: medida de dispersión, que indica qué tan dispersos están los datos con respecto a la media.

MODIFICACIÓN DEL PLANTEAMIENTO DE K-MEANS

En este caso, partimos del fichero de %avg de 168 líneas por nodo. Queremos establecer el nodo con sus $x=[168]$ posiciones, como elemento (ente a analizar), siendo cada una de las posiciones una característica (feature) del nodo.

La selección de características [10], puede resultar de mayor repercusión en el error de clasificación:

Pocas características, nos dan regiones de decisión más simples.

Las características deben constituir agrupaciones de valores medios que se mantienen muy semejantes entre objetos de la misma categoría y diferentes entre objetos de distintas categorías.

En nuestro caso son %average de cada hora, que coincide con esta definición.

Para este ensayo se va a usar el recurso [D.3]

Para cada línea, vamos a cotejar los resultados obtenidos, con la media, varianza, y desviación del nodo (que nos da un reflejo de su funcionamiento).

Aunque K-Means y PCA, parecen tener objetivos muy diferentes, hay una gran interconexión tal como se explica [5]

PCA busca representar todos los vectores de datos $n:n$ como combinaciones lineales de un pequeño número de vectores propios, y lo hace para minimizar el error de reconstrucción cuadrático medio.

K-Means busca representar todos los vectores de datos $n:n$ a través de un pequeño número de centroides de clúster, es decir, representarlos como combinaciones lineales de un pequeño número de vectores centroides de clúster .

PCA se utiliza para la reducción de dimensionalidad / selección de características, cuando el espacio de funciones contiene demasiadas funciones irrelevantes o redundantes o que no podemos utilizar. El objetivo es encontrar la dimensionalidad intrínseca de los datos.

Aquí hay un ejemplo bidimensional que se puede generalizar a espacios dimensionales superiores. El conjunto de datos tiene dos funciones, xx y yy , cada círculo es un punto de datos. Reducimos $v1$ y $v2$ a $v1$ por tener un mayor valor que $v2$ (realizando la reducción) y utilizando $v1$ como característica principal y reduciendo la dimensión y los valores a utilizar.

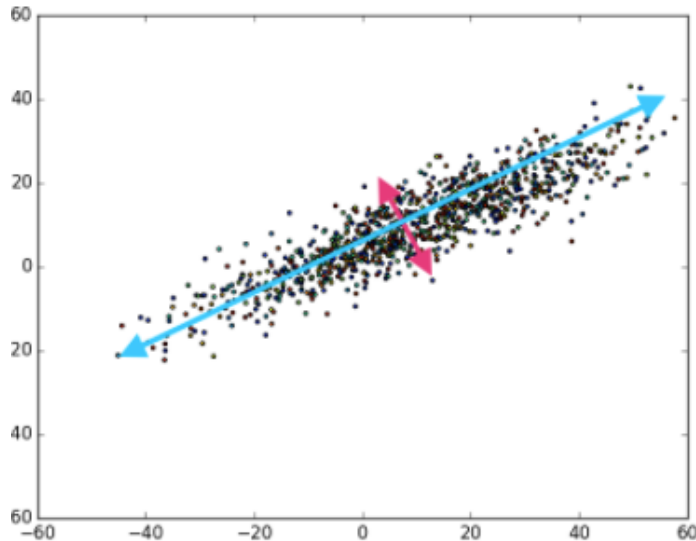


Figure 8 Ejemplo de reducción con PCA

Nuestras librerías de K-Means utilizan la reducción SVD (Single Valued Decompostion):

Obtenemos los siguientes valores para k (número de clusters)=3, 4 y 5:

nodo	media	varianza	Desviación	K=3	K=4	K=5
100019206	0.64	0.0378	0.1946	1	1	0
100021662	0.961	0.0005	0.0229	1	2	2
100052112	1.0	0.0002	0.0154	1	2	2
100055671	0.3515	0.0003	0.0173	2	0	1
100066658	0.781	0.0459	0.2144	1	1	0
100067629	0.0	0.0537	0.2317	0	3	4
100075173	0.308	0.0007	0.0267	2	0	1
100117943	0.0	0.0	0.0	0	3	4
100203236	0.2195	0.0056	0.0753	2	0	1
100205414	0.0915	0.0128	0.1133	2	3	4
100208668	0.3505	0.0048	0.0695	0	0	1
100224447	0.779	0.0009	0.0309	2	1	0
100233682	0.933	0.0051	0.0715	1	2	2
100252332	0.6235	0.0311	0.1765	1	1	0
100265331	1.0	0.0003	0.0199	1	2	2
100272294	0.478	0.0705	0.2656	2	1	0
100284574	0.978	0.0006	0.0251	1	2	2
100309685	0.24	0.0090	0.0953	2	0	1
100317715	1.0	0.0017	0.0412	1	2	2
100327214	0.0	0.0	0.0	0	3	4
100355392	0.689	0.0008	0.0294	1	1	0
100389873	0.121	0.0419	0.2048	2	0	1
100408768	0.0325	0.0151	0.1231	0	3	4
100413314	0.2155	0.0090	0.0952	2	0	1
100422568	0.0	0.0001	0.0128	0	3	4

100422926	1.0	0.2111	0.4594	1	1	3
100423635	0.167	0.0581	0.2410	2	0	1
100427449	0.0	0.0208	0.1443	0	3	4

Tabla 12 Listado de nodos, medias, varianza y desviaciones

K=3:

Clusters: [1 1 1 2 1 0 2 0 2 0 2 1 1 1 1 2 1 2 1 0 1 2 0 2 0 1 2 0]

Centroides

[[5.25357737 -0.05334234] [-4.69684157 -0.18139326] [2.1763397 0.28334617]]

Inertia⁸: 77.947

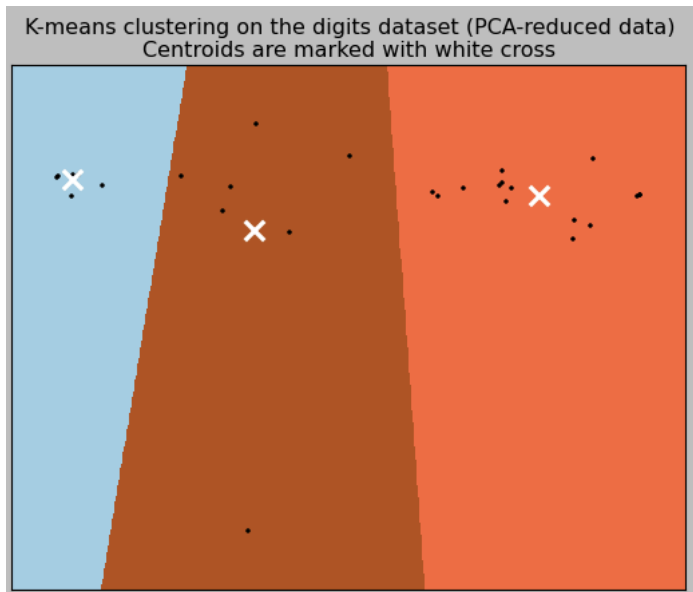


Figure 9 Estructura de los clusters y centroides k=3

K=4

Clusters: [1 2 2 0 1 3 0 3 0 3 0 1 2 1 2 1 2 0 2 3 1 0 3 0 3 1 0 3]

Centroides [[2.49906805 0.22231024] [-2.52370087 -0.50089586]

[-6.51694664 0.35019759] [5.25357737 -0.05334234]]

Inertia: 49.3379892552

⁸ Inertia: Cuanto menor es más coherente es el cluster. $\sum_{i=0}^n \min_{\mu \in C} (\|x_j - \mu_j\|)^2$

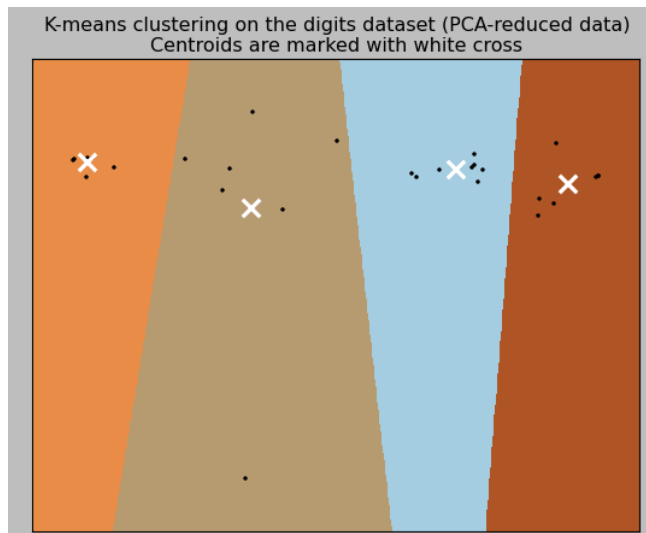


Figure 10 Estructura de clusters y centroides k=4

K=5

```
Clusters: [0 2 2 1 0 4 1 4 1 4 1 0 2 0 2 0 2 1 2 4 0 1 4 1 4 3 1 4]
Centroides [[-2.50133537  0.34774488] [ 2.49906805  0.22231024]
[-6.51694664  0.35019759] [-2.6578939  -5.59274034] [ 5.25357737  -
0.05334234]]
Inertia: 19.06895335
```

Vemos que:

- Clúster 4, es 0 o muy cerca de 0.
- Clúster 1, valores de media entre 0,1 y 0.3
- Clúster 2: Valores a 1 o muy cerca del 1.
- Clúster 3, valores intermedios entre el cluster 1 y el cluster 2.

El perfilado es correcto.

NUMERO DE CLUSTERS OPTIMO

Nos queda determinar cual debería ser el valor de k correcto. Aunque no exista un criterio objetivo para la selección del número de Clusters, si que se han implementado diferentes métodos que nos ayudan a elegir un número apropiado de Clusters para agrupar los datos; como son, el método del codo (elbow method), el criterio de Calinsky, el Affinity Propagation (AP), el Gap (también con su versión estadística), Dendrogramas,

Método del codo (Elbow Method)

Este método utiliza los valores de la inercia obtenidos tras aplicar el KMeans a diferente número de Clusters (desde 1 a N Clusters), siendo la inercia la suma de las distancias al cuadrado de cada objeto del Cluster a su centroide:

$$Inertia = \sum_{i=0}^N \|x_i - \mu\|^2$$

Una vez obtenidos los valores de la inercia tras aplicar el KMeans de 1 a N Clusters, representamos en una gráfica lineal la inercia respecto del número de Clusters. En esta gráfica se debería de apreciar un cambio brusco en la evolución de la inercia, teniendo la línea representada una forma similar a la de un brazo y su codo.

El punto en el que se observa ese cambio brusco en la inercia nos dirá el número óptimo de Clusters a seleccionar para ese data set; o dicho de otra manera: el punto que representaría al codo del brazo será el número óptimo de Clusters para ese data set.

De los datos obtenido de nuestros scripts obtenemos los valores de inercia:

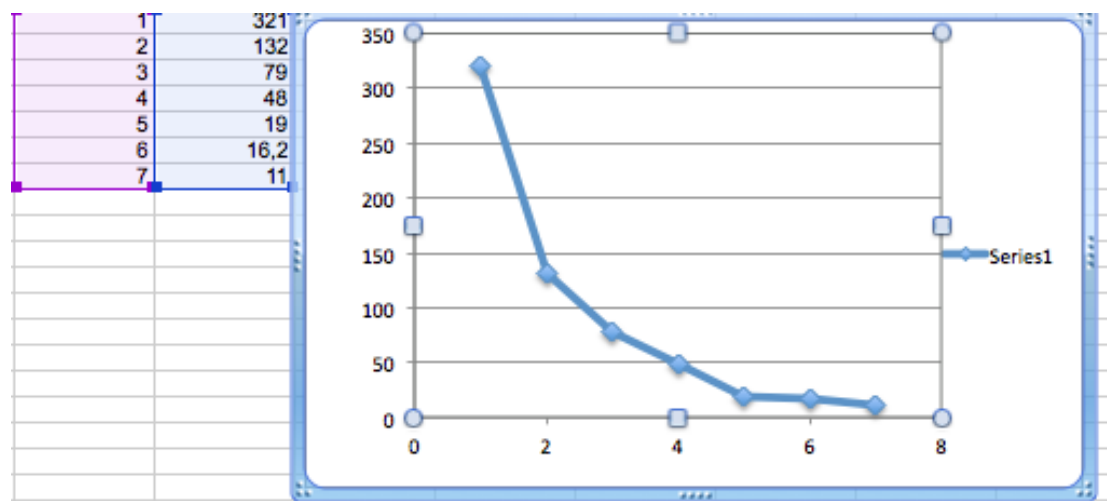


Figure 11 Obtención del número de nodos a través de la Inercia

El valor de clusters debería ser 3.

DB-SCAN

Partiendo de los mismos datos de partida, vamos a verificar cual es el número de cluster "k" adecuado utilizando [D.4] . Para lo cual vamos a usar el ente de Silhoutette⁹, (valor comprendido entre los rangos [-1, 1], que nos indica el correcto agrupamiento cuanto más cercano se encuentre de 1.

Valores de entrada :

- X: Vectores de %average.
- Desviación típica máxima de cada clúster: seleccionamos 0,15 en base a los datos de Tabla 12 Listado de nodos, medias, varianza y desviaciones.
- Eps: valor de la distancia máximo para encontrar al vecino=0,12
- Número de features=168
- Métrica para hallar la distancia: euclidean

Obteniendo los siguientes resultados:

Número óptimo de clúster: 4

V-measure: 0,967

Sihoutette coefficient: 0.709

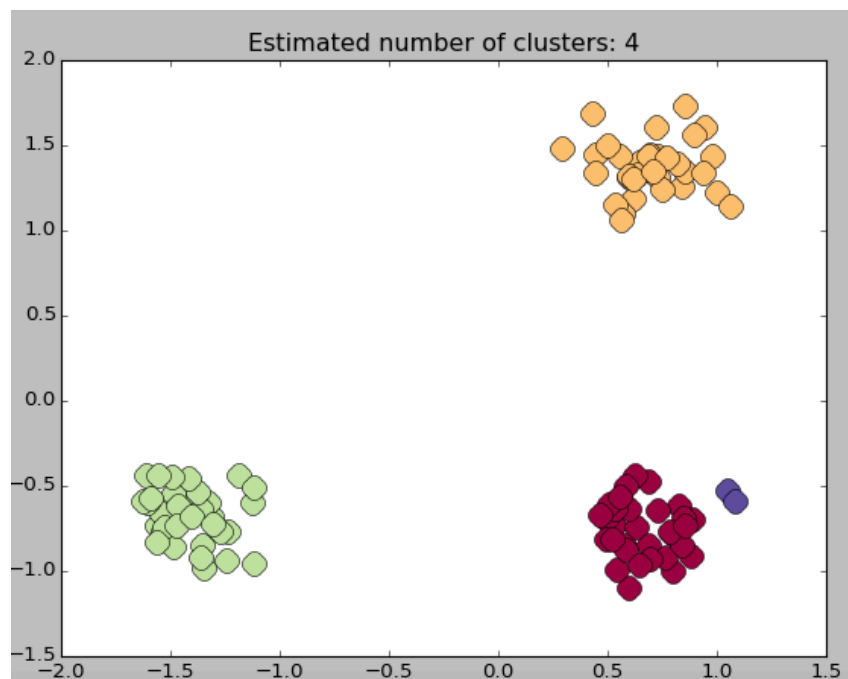


Figure 12 Db-scan clusters, eps=0,12

Encontrando el valor adecuado del coeficiente de Sihoutette en 0,2

V-measure: 1.000

⁹ Coeficiente de Silhoutette: $SC = \frac{1}{n} \sum_1^n s(x)$ con $s(x) = \frac{b(x)-a(x)}{\max\{a(x),b(x)\}}$

Silhouette Coefficient: 0.888
Estimated number of clusters: 3

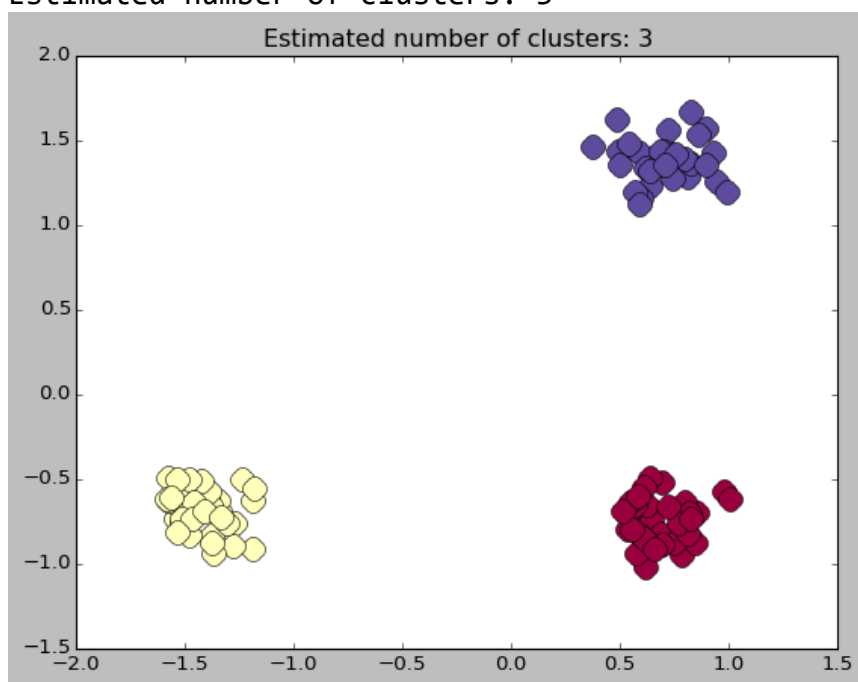


Figure 13 Db-scan número de cluster para eps=0,2

VERIFICACIÓN

VERIFICACIÓN PARA EL PROCEDIMIENTO DE PERFILADO GRUESO

1. Se pretende comprobar que uno o varios nodos se comporta similar a otro (distancia de Hamming pequeña o nula). Para lo cual, el se crearán plantillas de similitud, y se parsearán o pre-procesarán los nodos para disponer del fichero de vectores binarios de uptime con el formato adecuado y serán incluidos en el algoritmo.
De forma natural la distancia de Hamming de ambos será 0.

VERIFICACIÓN PARA EL PROCEDIMIENTO DE PERFILADO CON K-MEANS

1. Se pretende comprobar que uno o varios nodos se comporta similar a otro (que está comprendido dentro de un cluster. Para lo cual, el se crearán los clusters, y se parsearán o pre-procesarán los nodos para disponer del fichero de %average de uptime con el formato adecuado y serán incluidos en el algoritmo.
Se espera que el nuevo nodo se encuentre dentro del mismo cluster que el nodo referencia.

Nota: Al incluir varios nodos de un comportamiento distinto, la generación de los clusters (dinámicos) pueden variar. Esto no debiera ser problema por el alto volumen de datos de partida que procesamos.

Los resultados obtenidos han sido acordes al planteamiento, funcionando acorde a lo esperado.

Resultados:

Número de líneas de logs	Número de nodos comprendidos	Tiempos de resolución de K-MEANS (segundos)	Tiempo de resolución meta heurístico propio (segundos)
1.000	4	-	-
10.000	30	0,8	0,025
100.000	118	1	0.6
1.000.000	1.093	28	11

Tabla 13 Resumen de los tiempos de resolución del problema

Nota: Los datos del rendimiento del meta heurístico fueron descritos en Tabla 10 Comparaciones y tiempos de computación y el número de nodos en Tabla 4 Tamaños de las muestras disponibles

CONCLUSIONES

Para el caso de lo que hemos definido como “**procedimiento propio**”, su simplificación a un array de bit y su modelo comparativo:

- En primer lugar perfilado basado en su actividad (números de 1) que dispone el array de bits.
- En segunda instancia distancia comparativa a los demás nodos dentro de su perfil, hace que no le afecten los cambios de forma en la actividad de los nodos.

Se puede asemejar por su forma de trabajo a un cluster jerárquico.

Varía en que no agrupamos los nodos y comparamos el resultado con los demás. Esto nos hace tener comparativas de todos los elementos con todos los elementos, sin embargo para 10.000 líneas de logs procesadas en un sistema mono-hilo esta tarea tarda menos de 2 segundos.

La modificación y adecuación de esta tarea a un cluster jerárquico implicaría poco esfuerzo, sin embargo no creo que aportase mejora.

En el caso de la distribución de K-Means,:

Como se ha mostrado el comportamiento del algoritmo con K-Means utilizando todo el vector (siendo cada hora una característica del nodo), y realizando la reducción a través de PCA. nos permite identificar los cluster y perfilar un nuevo nodo dentro del mismo tras identificar su comportamiento, Su rendimiento es adecuado perfilando 10.000 líneas de los logs en menos de 2 segundos.

Por lo cual se puede usar la solución con un mes (array de 512 elementos) siempre que se definan como características.

Este planteamiento puede extenderse a otras soluciones de perfilado que utilicen en los logs el timestamp en formato epoch como por ejemplo los sistemas de Single Sing On, pudiendo usar esta metodología para identificar el comportamiento de un mismo usuario y establecer módulos de autorización basados en el comportamiento.

PROXIMOS PASOS:

Se plantean dos ámbitos de mejora del análisis:

En primer lugar *analizar si modelos mixtos pueden llegar mejora la clasificación* al elegir el elemento central para el cálculo de distancias, a priori se estima que no para valores aleatorios.

Analizar si otro el *modelo probabilístico dinámico* EM tal como se ha planteado en 15, es aplicable. A priori creemos que puede ser una buena alternativa.

En segundo lugar, *mejorar la gestión de acceso a los logs*:

En este apartado se pretende dar una alternativa a las limitaciones encontradas de cara a poder procesar ingentes cantidades de datos en la capacidad de almacenar la información, la calidad y estandarización de los procedimientos y velocidad que nos aportan los recursos u orígenes de la información.

Históricamente en los sistemas los logs han almacenado en texto plano, y posteriormente de forma histórica en BBDD por su disposición a generar datos indexados y vistas con la información pre-procesada. Realizando esta tarea a nivel de recurso y eliminando esta complejidad en el desarrollo de los algoritmos.

Seti@home nos propone estas dos características (18GB en texto plano y 21 GB para la BBDD de mysql).

Ha habido avances de centralización de los mismos y recolección de estas evidencias en entornos distribuidos y con distintos orígenes de información (mezcla de logs en texto plano, BBDD) como la herramienta de graylog e incluso otras opciones como la inclusión de datos no estructurados.

Sin embargo la mayor limitación que podemos encontrar son las consultas a la información.

Tras la incursión de las bases de datos noSQL, se eliminan parte de las deficiencias y limitaciones en los que algunos algoritmos estaban etiquetados y que se han analizado en este estudio.

Hemos de definir sus características de las bases de datos noSQL, ya que a diferencia de las SQL, no todas comparten el mismo lenguaje y no todas trabajan del mismo modo.

Las hay que muestran los datos en un grafo y las hay que lo hacen en documento.

Las bases de datos NoSQL que más se utilizan en entornos productivos de gran carga en procesamiento son:

Tabla 14 Beneficio y contras en las distintas versiones de BBDD noSQL

Nombre	Descripción	Pros	Contras
Neo4J	Base de datos que muestra los datos en un grafo.	Rapidez en las búsquedas (de las más rápidas). Lenguaje sencillo (Cipher, el más intuitivo y fácil de	

		aprender). App web intuitiva. Api sencilla de utilizar. Licencia GNU (Copyleft).	
OrientDB	Base de datos que muestra los datos en un grafo, forma documento y key value (la más completa de todas)	Rapidez en las búsquedas. Lenguaje. sencillo (Cipher) App web intuitiva. API muy sencilla de utilizar.	Consola web poco consistente. Distinción entre versiones community y profesional.
MongoDB	Base de datos que muestra los datos en documento	Rapidez en las búsquedas. API sencilla de utilizar.	Las búsquedas se ven condicionadas por el modo en que se almacenan la información en la BBDD.. Restringida a licencia.
Cassandra	Base de datos key value	Rapidez en las búsquedas. Licencia Apache Software Foundation	Configuración

COMPARATIVA DE BUSQUEDA SQL VS noSQL

En la siguiente tabla se muestran los tiempos de respuesta entre realizar una búsqueda en una BBDD relacional y una no relacional (Neo4J en particular):

Nivel Profund	Relacional Tiempo (seg)	Neo4j Tiempo (seg)	Registros Devueltos
2	0.016	0.01	≈ 2.500
3	30.267	0.168	≈ 110.000
4	1543.505	1.359	≈ 600.000
5	No Termina	2.132	≈ 800.000

*Partner y Vukotic's Experiment - Graph Databases by Ian Robinson, Jim Webber, and Emil Eifrem

Figure 14 Comparativa de rendimiento SQL NoSQL

Como se puede apreciar, con un número pequeño de búsquedas los tiempos son casi idénticos, pero al realizar una búsqueda que requiera un gran volumen de datos a devolver, los tiempos son completamente diferentes.

BASES DE LA SELECCIÓN DE ESTA SOLUCIÓN

La solución debe disponer de las siguientes características:

- Tiempo de respuestas altos.
- Use una lógica que permita realizar cambios en el código sin dificultad.
- Disponga una API que permita un uso sencillo.

- Limitación de dependencias. Por ejemplo, al eliminar un dato en una solución SQL, ha de eliminarse de todas las tablas donde se pueda encontrar y/o las referencias. Las bases de datos NoSQL no tienen esta limitación, ya que no existe el join.

Se propone una base de datos NoSQL de Neo4J (como alternativa de almacenar los logs), por no disponer de las limitaciones anteriormente descritas, disponer de un motor de persistencia y una consola web amigable para su gestión. Permite usar una lógica a nivel de arquitectura sencilla, dispone de una API de fácil integración con tiempos de aprendizaje cortos. El lenguaje de consulta ciper, es muy extendido y sencillo de aprender y utilizar.

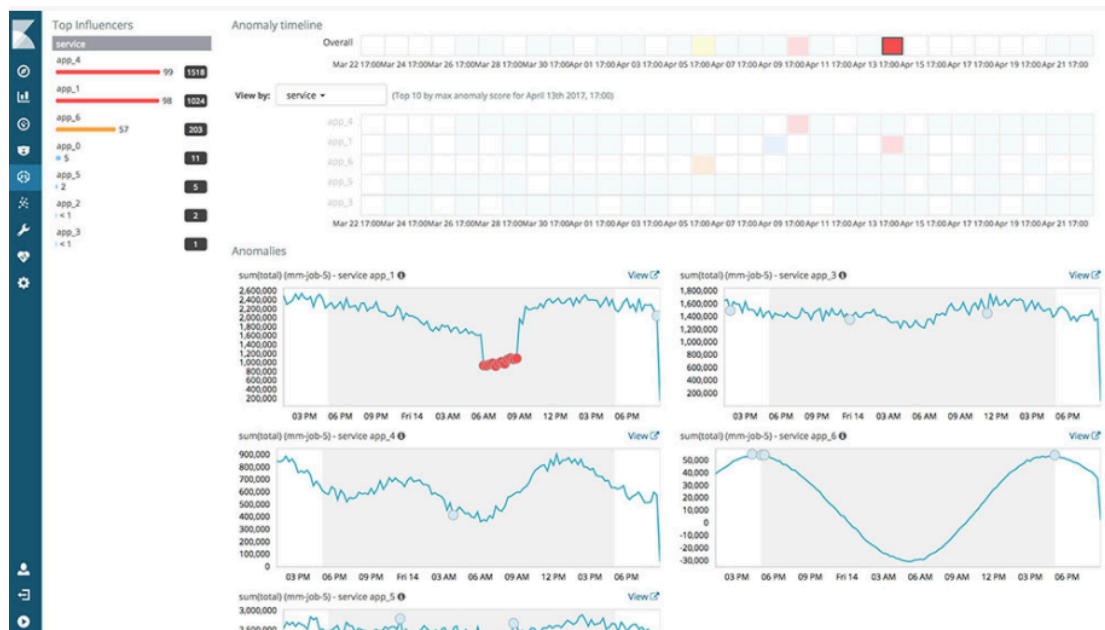
Dispone de un licenciamiento AGPLv3 (GNU) que está dentro del marco de Copyleft.

Su última versión estable es la 3.3.0 de Octubre del 2017.

SOLUCIONES COMERCIALES QUE APOYAN LA TENDENCIA

Herramientas comerciales para la clusterización o perfilado de datos masivos: <https://www.elastic.co>

Dispone de varias soluciones o componentes, que unidos forman una solución de matching Learning para la gestión de patrones tanto de actividad de las aplicaciones, como el perfilado de tendencias.



Componentes:

- Elasticsearch (motor de búsqueda)
- Kibana (Dashboard)
- Beats (agentes para el paso de datos a Elasticsearch o Logstash)

POSIBLE IMPLEMENTACIÓN PROPIA

Estas tecnologías se basan en poder definir nuestra propia solución, basándonos en los componentes que utilizan las herramientas comerciales, y asociadas a la explotación masiva de logs para millones de entradas_

Componente Elasticsearch:

ElasticSearch es un motor de búsqueda, que sirven para indexar entornos distribuidos, orientado para grandes volúmenes de datos (escalabilidad horizontal y vertical) y que nos ofrece un acceso en tiempo real y orientado a documentos (JSON) y que no trabaja con schemas de base de datos (estructuras fijas) buscando por patrones. Será un complemento para la solución.

Componente Logstash:

Logstash, herramienta que funciona bajo un JVM que nos permitirá administrar los logs. Nos aporta la posibilidad de usar codecs (permiten cambiar el formato de entrada a otro formato, por ejemplo los datos no son de texto o cambio de estructura). Dispone conectores de entrada y de salida (para almacenar los logs parseado por los codecs). Y filtros (programados en Grok, que nos ayuda a parsear o filtrar las entradas.

Componente Kibana:

Es un dashboard que nos facilitará la gestión de la información de elastic search.

Neo4j se integra directamente con Logstash, elasticsearch tal como indica su propia web, así como con Python, para poder usar las utilidades de stickit.

GLOSARIO

B

BOINC: (Berkeley Open Infrastructure for Network Computing) es un sistema open source liberado bajo la licencia LGPL que permite la computación en grilla voluntaria y permite a los investigadores acceder a abundantes recursos de procesamiento obtenidos de la suma de dispositivos personales en todo el mundo.

C

CF: Cluster Features, o características del cluster.

COMPUTACION EN GRID: es un sistema para la compartición de recursos de procesamiento y almacenamiento de varios servidores con un alto rendimiento y bajo demanda. No se trata de varios equipos en paralelo, sino de un equipo virtual que va balanceando recursos libres de muchos servidores para ofrecernos un conjunto escalable de recursos conforme van siendo necesarios. De esa manera, aprovechando recursos disponibles de entre los recursos ofertados, para montar una potente infraestructura virtual distribuida en varios equipos físicos. El concepto de grid no se debe asociar a computación en la nube.

D

DBSCAN: Density Based Spatial Cluster Applications with Noise (KDD 1996)

DHT: (Distributed Hash Table) son un tipo de tablas de hash, almacenan pares (clave, valor) y permiten consultar el valor asociado a una clave, en las que los datos se almacenan de forma distribuida en una serie de nodos (sistemas distribuidos) y proveen un servicio eficiente de búsqueda que permite encontrar el valor asociado a una clave. Para esto último usan un sistema de enrutado que permite encontrar de forma eficiente el nodo en el cual está almacenada la información que se necesita.

DENCURE: Density based Clustering (Hinnedburg & Keim, KDD 1998).

F

FAR: Falsa aceptación.

FDP: Función de densidad probabilística.

FRR : Tasas de Falso rechazo.

O

OPTICS: Ordering Points To Identify The Cluster Structure (Ankerst et al SIGMOD 1999)

P

PCA: Principal Component Analysis

PODS: son servidores autónomos en los que los usuarios pueden registrarse y comunicarse con toda la red de microblogging.

AGRADECIMIENTOS

Inicialmente a mi tutor que ha coordinado con paciencia e ilusión la investigación.

En segundo lugar a la institución de la UOC que potencia el desarrollo de aplicaciones propias y permite a los estudiantes basar sus investigaciones que aporten al desarrollo de otros compañeros.

REFERENCIAS

Referencias base:

- [1] Briony J Oates: Researching Information Systems and Computing, 2006, ISBN 978-1-4129-0223-6

Otras referencias:

- [2] Charu C. Aggarwal & Chandan K. Reddy (editors): Data Clustering: Algorithms and Applications. Chapman & Hall / CRC Press, 2014. ISBN 1466558210
- [3] Charu C. Aggarwal Aggarwal & Jiawei Han(editors editors): Frequent Frequent Pattern Pattern Mining. Springer Springer, 2014. , 2014. ISBN 3319078208.
- [4] Ding, C. He. X. [2004], K-Means Clustering via Principal Component Analysis.
- [5] Fisher, D., Knowledge acquisition via incremental conceptual clustering, Machine Learning, VOL. 2, 139-172, 1987.
- [6] Garre, M., Cuadrado, J, Sicilia, M.A, [artículo].“Comparación de diferentes algoritmos de clustering en la estimación de coste en el desarrollo de software”.
- [7] Gutiérrez R., González A., Torres, F. J.A. Gallardo (1994)., [artículo], “Técnicas de Análisis de datos Multivariable. Tratamiento computacional”. Universidad de Granada.
- [8] Jain A.K., Murty M. N., and Flynn P. J. Data clustering: A review. ACM Computing Survey, 31(3):264–323, 1999.
- [9] Lloyd S.P.: Least squares quantization in PCM. Bell Telephone Laboratories, 1957. Published much later in IEEE Transactions on Information Theory 28 (2): 129–137, 1982. DOI 10.1109/TIT .1982.1056489

- [10] MacQueen J. B.: Some Methods for classification and Analysis of Multivariate Observations. Proceedings of 5th Berkeley Symposium on Mathematical Statistics and Probability 1, pp. 281–297. University of California Press, 1967. MR 0214227. Zbl 0214.46201
- [11] Villagra A., Guzmán A. , Pandolfi D. , [artículo], “Análisis de medidas no-supervisadas de calidad en clusters obtenidos por K-means y Particle Swarm Optimization”.Universidad Nacional de San Luis.
- [12] Lazaro D., [tesis], “A Middleware for Service Deployment in Contributory Computing System”. 2004

Referencia gris:

[13] Gómez, H.,[documento interno], “Technical Report of the Garlanet profiling process”.

RECURSOS:

[A] Setti@Home, Website. (<https://setiathome.berkeley.edu/>)

[B] BOINC, Website. (<https://juncotic.com/boinc-computacion-distribuida-voluntaria/>)

[C] Scikit Learn: <http://scikit-learn.org>

[D] Recursos propios:

[D.1] Pre-procesado de la información:

De esta tarea se encargará un conjunto de scripts en Perl:

1-array_week.pl:

Su función es leer el fichero de logs y generar :

- un nuevo fichero con los datos de actividad de cada nodos, y con un formato de (Año, mes, día, hora, minuto) de actividad (parada-comienzo), se extrae una semana completa de actividad. Comenzando el mismo día de la semana extraído por los logs. Generará el fichero salida_vector_avg.txt (destinado al análisis de K-Means).
- Un fichero similar pero con el ajuste binario, creando el fichero de salida vector_array_bist.txt (para el procedimiento de perfilado grueso).

Nota: Esta adecuación se basa en que podríamos modificar el algoritmo para el análisis de un día, o mes, comenzando siempre el análisis un Lunes.

Esto nos alteraría el tamaño de bits a almacenar en los arrays. Este procedimiento se ha seguido para proteger líneas repetidas en el fichero origen de los logs.

2-correlate-avg.pl:

Su función es prevenir líneas duplicadas y eliminar los outliers definidos como aquellos nodos sin actividad o siempre activos.

Adicionalmente:

- Agrega los promedios de los nodos que tengan más de una entrada en la misma hora.
- Ordena y elimina los nuevos outliers tras acumular los nuevos promedios de actividad para los ficheros de %average y vector de bits.

Genera varios ficheros de trace y los ficheros:

- Outpunt_avg_wihout_outliers.txt
- Vector_bits_without_outliers.txt
- Outliers.txt (que nos lista los host que no vamos a procesar).

3-prepare-for-kmeans.pl:

Su función es adecuar el fichero con el vector de 168 posiciones con el %average que generará un fichero con el formato requerido para el procedimiento de K-Means.

- Nos muestra los valores de media, varianza y desviación estándar de cada vector de bits y el número de ellos que se aleja de los límites superiores e inferiores.
- Nos aporta, la media, varianza y desviación estándar del conjunto total de datos del análisis.

5-density.pl:

Nos calcula la densidad de los arrays.

[D.2] 4-hamming-distance.pl:

Script que da los resultados del procedimiento propio.

[D.3] k-meanas-feature_vector.py

Script de Python haciendo uso de la librería de skearn (skearn.cluster) para importar el módulo ya creado de K-Means. También usaremos sklearn.decompostion para importar el módulo de PCA.

Nota: El uso de las librerías aquí mencionadas mejora drásticamente el rendimiento agrupando los clusters en menos de 1 segundo.

[D.4] db-scan.py

Script Python calcular los clusters basándonos en le agrupamiento por densidad DB-SCAN 10.