

Plan de Proyecto de Titulación

Cristian Pachacama

May 25, 2018

1 INFORMACIÓN BÁSICA

1. **Propuesto por:** Cristian David Pachacama Simbaña
2. **Línea de Investigación:** Estadística Aplicada
3. **Fecha:** 11 de mayo de 2018

2 RELACIÓN

1. **Nombre del Proyecto de Investigación:** Forecast and Impact of extreme low levels of streamflow in hydropower plants.
2. **Director del Proyecto de Investigación:** PhD. Adriana Uquillas

3 INFORMACIÓN DEL TRABAJO DE TITULACIÓN

3.1 Título del Trabajo de Titulación

ANÁLISIS CLUSTER DE SERIES DE TIEMPO.

3.2 Planteamiento del Problema

Brasil tiene una de los sistemas hidrológicos más complejos, diversos y extensos del mundo. A diferencia de la gran mayoría de los países desarrollados, Brasil tiene en los ríos su principal fuente de generación de electricidad, ocupando el tercer lugar dentro de los más grandes productores hidroeléctricos del mundo. Debido a la importancia del sector hidroeléctrico, buscar formas de facilitar y mejorar el modelamiento de datos asociados a este sector es un problema prioritario. Provocado por la dificultad que supone lidiar con la enorme cantidad de datos asociados a mediciones de Caudales de los ríos que componen este sistema, que cuenta con alrededor de 150 estaciones de medición repartidas en todo Brasil. Dichos datos se presentan en forma de Series de Tiempo que posee tres características que dificultan su análisis, la primera es que estas series de tiempo poseen observaciones diarias de los caudales en un periodo de tiempo de alrededor de 30 años, es decir, son series muy extensas. La segunda característica es que estas series de tiempo son estacionales, y por último existe evidencia de que el ruido o error asociado a estas series no se distribuye normalmente sino que su distribución posee colas más pesadas como las analizadas en teoría de valores extremos.

En ese contexto, notamos que es posible disminuir la dimensión del problema a través la identificación de clústers o zonas representativas (no necesariamente geográficas) que resuman comportamiento temporal que poseen los caudales de los ríos.

3.3 Justificación

Ya que el problema se basa en identificar grupos de ríos cuyos Caudales se comportan de manera similar en el tiempo, se propone la utilización de el “Análisis Clúster de Series de Tiempo”, que es una técnica de agrupamiento que considera una función de “disimilitud” entre las series de tiempo (que mide que tan distintas son un par de series) y a partir de ella crea grupos de series, cada grupo contiene series de tiempo “parecidas”. Al elegir adecuadamente la función de disimilitud (diseñada para series de tiempo) es posible agrupar a los ríos en grupos basados en el comportamiento temporal de sus caudales. Esto con la finalidad de lidiar con la complejidad que supone analizar la enorme cantidad de series de tiempo de caudales, pasando de analizar alrededor de 150 series a unas pocas (una serie por Clúster), sin dejar de lado la estructura y comportamiento estacional de cada una de ellas, partiendo de una adecuada elección de la función de disimilitud.

3.4 Hipótesis

3.5 Objetivo General

El principal objetivo del proyecto es el de agrupar estaciones (asociadas a ríos en Brasil) en zonas, basandonos en el comportamiento temporal del caudal de los ríos que se mide diariamente en dichas estaciones, a fin de facilitar el modelamiento de los caudales de alrededor de 150 estaciones repartidas en todo Brasil, que como sabemos posee uno de los sistemas hídricos más grandes y complejos del mundo, por lo que al estudiar su comportamiento es posible escalarlo y adaptarlo a sistemas menos complejos como el de nuestro país.

3.6 Objetivos Específicos

1. Comparar una gama de técnicas tanto de clusterización, así como la elección de distintas métricas o funciones de disimilitud a fin de identificar aquella combinación (disimilitud/algoritmo de clusterización) que permita un adecuado agrupamiento de las series de tiempo asociadas al Caudal de los ríos, que tienen la peculiaridad de que las series están conformadas por observaciones diarias de los caudales y además suelen ser series estacionales, y este tipo de series no han sido analizadas con esta técnica previamente.
2. Comparar el modelamiento de estas series de tiempo usando el análisis cluster versus el modelamiento sin un previo análisis, a fin de mostrar las ventajas en cuanto a tiempo y eficiencia de la metodología planteada.
3. Automatizar la descomposición y Análisis Clúster de series de tiempo de Caudales a través de la utilización de software estadístico y programación.
4. Crear una Plataforma de Análisis (Visualización y descomposición) de las Series de Tiempo de Clima, Caudales, y Contaminación, que permita a más investigadores analizar toda la información recolectada durante este proyecto.

3.7 Metodología

El Análisis Clúster es una técnica de aprendizaje no supervisada que tiene como objetivo dividir un conjunto de objetos en grupos homogéneos (clústers). La partición se realiza de tal manera que los objetos en el mismo clúster son más similares entre sí que los objetos en diferentes grupos según un criterio definido. En muchas aplicaciones reales, el análisis de clúster debe realizarse con datos asociados a series de tiempo. De hecho, los problemas de agrupamiento de series de tiempo surgen de manera natural en una amplia variedad de campos, incluyendo economía, finanzas, medicina, ecología, estudios ambientales, ingeniería y muchos otros. Con frecuencia, la agrupación de series de tiempo desempeña un papel central en el problema estudiado. Estos argumentos motivan el creciente interés en la literatura sobre la agrupación de series de tiempo, especialmente en las últimas dos décadas,

donde se ha proporcionado una gran cantidad de contribuciones sobre este tema. En [?] se puede encontrar un excelente estudio sobre la agrupación de series de tiempo, aunque posteriormente se han realizado nuevas contribuciones significativas. Particularmente importante en la última década ha sido la explosión de documentos sobre el tema provenientes tanto de comunidades de minería de datos como de reconocimiento de patrones. [?] proporciona una visión general completa y exhaustiva de las últimas orientaciones de minería de datos de series de tiempo, incluida una gama de problemas clave como representación, indexación y segmentación de series de tiempo, medidas de disimilitud, procedimientos de agrupamiento y herramientas de visualización.

Una pregunta crucial en el Análisis Clúster es establecer lo que queremos decir con objetos de datos "similares", es decir, determinar una medida de similitud (o disimilitud) adecuada entre dos objetos. En el contexto específico de los datos asociados a series de tiempo, el concepto de disimilitud es particularmente complejo debido al carácter dinámico de la serie. Las diferencias generalmente consideradas en la agrupación convencional no podrían funcionar adecuadamente con los datos dependientes del tiempo porque ignoran la relación de interdependencia entre los valores.

De esta manera, diferentes enfoques para definir una función de disimilitud entre series de tiempo han sido propuestos en la literatura pero nos centraremos en aquellas medidas asociadas a la autocorrelación (simple, e inversa), correlación cruzada y periodograma de las series (Ver: [?]; [?]; [?]; [?]). Estos enfoques basados en características tienen como objetivo representar la estructura dinámica de cada serie mediante un vector de características de menor dimensión, lo que permite una reducción de dimensionalidad (las series temporales son esencialmente datos de alta dimensionalidad) y un ahorro significativo en el tiempo de cálculo, además de que nos ayudan a alcanzar el objetivo central por el que usaremos el Análisis Clúster que es el de la modelización de series de tiempo.

Una vez que se determina la medida de disimilitud, se obtiene una matriz de disimilitud inicial (que contiene la disimilitud entre parejas de series), y luego se usa un algoritmo de agrupamiento convencional para formar los clústers (grupos) con las series. De hecho, la mayoría de los enfoques de agrupamiento de series de tiempo revisados por [?] son variaciones de procedimientos generales como por ejemplo: K-Means, K-Medoids, PAM, CLARA [?] o de Clúster jerárquico que utilizan una gama de disimilitudes específicamente diseñadas para tratar con series de tiempo y algunas de sus características.

Una etapa adicional dentro del análisis consiste en determinar la cantidad de clústers que es más apropiada para los datos. Idealmente, los clústers resultantes no solo deberían tener buenas propiedades estadísticas (compactas, bien separadas, conectadas y estables), sino también resultados relevantes. Se han propuesto una variedad de medidas y métodos para validar los resultados de un análisis clúster y determinar tanto el número de clústers, así como identificar qué algoritmo de agrupamiento ofrece el mejor rendimiento, algunas de estas ellas pueden encontrarse en [?]; [?] ; [?] ; [?]. Esta validación puede basarse únicamente en las propiedades internas de los datos o en alguna referencia externa.

3.8 Plan de Trabajo

El plan de trabajo se resume en la Tabla ??.

Actividad	
Obtención de Datos	Descarga de Bases de Datos Brasil de las fuentes:Datos Hídricos Operado
Tratamiento de Datos	Depuración de Datos. Estandarización de Datos. Unificación de Datos par
Análisis y Modelamiento	Análisis y comparación de métricas y funciones de disimilitud. Comparació
Automatización	
Conclusiones y Recomendaciones	

Table 1: Plan de Trabajo

3.9 Cronograma

A continuación se muestra el cronograma de actividades planificadas resumidas en la Tabla ??

	Abril	Mayo	Junio	Julio	Agosto
Obtención de Datos	X	X			
Tratamiento de Datos		X			
Análisis y Modelamiento		X	X	X	
Automatización			X	X	X
Conclusiones y Recomendaciones					X

Table 2: Cronograma de Trabajo