

# IDENTIFICAÇÃO DE GRUPOS REGIONAIS DE CHUVA POR ANÁLISE DE AGRUPAMENTOS

Frederico Ivanchechen de Mattos<sup>1</sup>, Mino Viana Sorribas<sup>2</sup>, Einara Zahn<sup>2</sup>  
e José Eduardo Gonçalves<sup>2</sup>

<sup>1</sup>Programa de Pós-Graduação em Engenharia de Recursos Hídricos e Ambiental (PPGERHA),  
Universidade Federal do Paraná (UFPR), Brasil.

<sup>2</sup>Sistema Meteorológico do Paraná (SIMEPAR), Brasil.

E-mail: fivanche@gmail.com, mino.sorribas@simepar.br, einara.zahn@gmail.com, jose.eduardo@simepar.br

## Introdução

Medições manuais e automáticas de precipitação concedem informações cruciais para o gerenciamento de recursos hídricos e projetos de obras hidráulicas. Entretanto, problemas na telemetria, entupimento dos pluviômetros e avarias dos instrumentos de medição durante tempestades podem causar falhas e/ou erros na série de dados. Consequentemente, a Agência Nacional de Águas (ANA) exige a consistência desses dados em antemão a sua disponibilização. Nesse sentido, métodos de análise de frequência regional são utilizados para identificar erros e, na escala temporal adequada, preencher falhas. O requerimento primário para essa análise é a identificação de áreas pluviometricamente homogêneas. Nesse estudo de caso explora-se a análise de agrupamentos pelo método de Ward utilizando a distância Euclidiana como medida de dissimilaridade e parâmetros de validação de agrupamentos para determinar o número ótimo de grupos. Apesar da forte semelhança entre as observações de todos os postos, a análise de agrupamento isolou sete grupos regionais, reduzindo a arbitrariedade na tomada de decisões.

Palavras-chave: Análise de grupos, Método de Ward, Hidrologia estatística.

## Metodologia

### Dados de entrada

Séries de precipitação acumulada mensal foram levantadas entre 2000 e 2015 para 30 estações pluviométricas monitoradas pela CTG-Brasil (Chinese Three Gorges Corporation) e mais uma adicional pertencente à rede do SIMEPAR (Sistema Meteorológico do Paraná). Todas as estações se encontram na bacia hidrográfica do Paranapanema, um rio de 929 km de extensão que divide os estados de São Paulo e Paraná, drenando cerca de 100800 km<sup>2</sup> (ANA, 2016). Para cada posto  $x_i$  foram acumulados os registros diários de precipitação no mês  $j$ , obtendo-se o valor  $z_{ij}$ . Como a análise de agrupamentos exige que haja observações no mesmo mês para todas as estações, 5% dos dados foram invalidados devido a falhas no histórico.

### Análise de agrupamentos

Conquanto existam diversos métodos de identificação de regiões homogêneas (Satyanarayana e Srinivas, 2008), a análise de agrupamentos é amplamente empregada para dados de precipitação em estudos no mundo inteiro (Pelczer e Cisneros-Iturbe, 2008; Santos, Lucio e Silva, 2015; Shirin e Thomas, 2016; Terassi e Galvani, 2017). Há duas vertentes da análise, ambas com o mesmo objetivo: formar grupos cujos membros pertencentes sejam os mais semelhantes possíveis entre si e os mais distintos dos demais. Métodos hierárquicos, como o de Ward, iniciam com cada estação em um grupo (de tamanho um, denominado singleton) e, à medida que se itera, o método funde os grupos mais semelhantes. Essa união pode ser baseada em uma gama de princípios, como mínima distância entre grupos (single linkage) e menor distância máxima entre grupos (complete linkage). A grande vantagem do método de Ward é

que essa junção se baseia no mínimo incremento de variância, portanto, procurando minimizar a distância média entre as séries históricas dos postos. Foi adotado a fórmula de Lance-Williams para o algoritmo de Ward, conforme a equação [1]:

$$d(C_i \cup C_j, C_k) = \frac{(|C_k| + |C_i|)d(C_k, C_i)^2 + (|C_k| + |C_j|)d(C_k, C_j)^2 + |C_k|d(C_i, C_j)^2}{|C_i| + |C_j| + |C_k|} \quad [1]$$

Em que  $d(C_i \cup C_j, C_k)$  representa o custo de fusão do grupo  $C_i$  ao  $C_j$  em relação ao grupo  $C_k$ , que não pertence a nenhum dos outros dois.  $|C_s|$  representa o número de membros no grupo, enquanto que  $d(C_a, C_b)$ , a distância entre os grupos  $C_a$  e  $C_b$ .

### Medida de Dissimilaridade

A distância, ou dissimilaridade, entre grupos pode ser medida por métricas como distância Euclidiana (Santos, Lucio e Silva, 2015), distância Euclidiana Padronizada (Satyanarayana e Srinivas, 2008), distância Mahalanobis (Detzel, Bessa e Mine, 2013) e outras. Optou-se pela distância Euclidiana pelo fato dela poder ser entendida como uma medida em um espaço multidimensional (Terassi e Galvani, 2017) e ter sido amplamente utilizada em estudos similares. Sua fórmula é expressa em [2].

$$d(x_i, x_j) = \sqrt{\sum_{k=1}^{n_j} (z_{ik} - z_{jk})^2} \quad [2]$$

Onde  $x_i$  representa o vetor de observações  $(z_{i1}, z_{i2}, \dots, z_{in_j})$ .

### Determinação do número ótimo de grupos

A inspeção visual de dendrogramas, como o da figura 1, é o método mais simples para determinar o número de grupos. Ele é um gráfico em formato de árvore cujas folhas representam os postos pluviométricos e a união dos galhos, os pontos de mesclagem dos grupos. Espera-se que a distância vertical entre os nódulos seja pequena na base do gráfico e que haja um salto logo após atingir-se o número ótimo de agrupamentos. Para reduzir a incerteza, recorreu-se a três parâmetros de desempenho de agrupamentos sugeridos pela literatura técnica (Pelczer e Cisneros-Iturbe, 2008; Satyanarayana e Srinivas, 2008), todos partindo da premissa que bons agrupamentos são compactos e distantes entre si. O índice Dunn (1974) é a razão entre a menor distância entre grupos e a maior distância entre membros de um mesmo grupo, sendo calculado pelas equações [3], [4] e [5]. Quanto maior o índice de Dunn (DI), melhor o agrupamento. O índice de silhueta (SI) também relaciona distância intergrupos e intragrupos, embora com as equações [6], [7] e [8]. SI Próximos da unidade apontam bons agrupamentos, os próximos de zero indicam grupos sobrepostos e valores negativos são sinais de alocação errônea de postos nos grupos. Satyanarayana e Srinivas (2008) usam o índice Davies-Bouldin (DBI), uma razão entre desvios das observações dos membros em relação ao centróide de seu grupo e a distância entre grupos. Verificando sua formulação nas equações [9] e

