

Exercício Programa 4

Ana Paula Lopes Cavalcante
Gabriel Tetsuo Haga

Junho 2021

1 Introdução

O objetivo deste relatório é apresentar o programa feito para solucionar o seguinte problema: Considere um modelo estatístico m-dimensional multinomial, com observações x , e informações *a priori* y e parâmetro θ , tal que $x, y \in \mathcal{N}^m$ e $\theta \in \Theta = S_m = \{\theta \in R_m^+ | \theta'1 = 1\}$, $m = 3$. E esse modelo segue a função densidade de probabilidade *a posteriori* de Dirichlet dada por:

$$f(\theta|x, y) = \frac{1}{B(x+y)} \prod_{i=1}^m \theta_i^{x_i+y_i},$$

sendo $B(x+y)$ a função Beta multivariável. Considere agora duas funções uma $T(v)$ e $W(v)$, definidas como:

$$T(v) = \{\theta \in \Theta | f(\theta|x, y) \leq v\}$$

$$W(v) = \int_{T(v)} f(\theta|x, y) d\theta$$

Assim o problema proposto consiste em criar uma função $U(v)$ que estime o valor de $W(v)$. Para tal deve-se definir k pontos de cortes: $0 = v_0 < v_1 < v_2 < \dots < v_k = \sup f(\theta)$, e também deve-se criar n pontos θ_i , sendo $i \in \{1, 2, 3, \dots, n\}$ que siga a distribuição de Dirichlet. Após isso estima-se que o $W(v_j) - W(v_{j-1})$ será, aproximadamente, a fração $p = \frac{l}{n}$, tal que $l = \sum_{i=1}^n R(\theta_i, v_{j-1}, v_j)$, sendo a função $R(\theta, v_{j-1}, v_j)$ definida como:

$$R(\theta, v_{j-1}, v_j) = \begin{cases} 1, & \text{se } v_{j-1} \leq f(\theta) \leq v_j \\ 0, & \text{se } f(\theta) < v_{j-1} \text{ ou } f(\theta) > v_j \end{cases}$$

A função deverá também ajustar, dinamicamente, os cortes v_j , para tal que $W(v_j) - W(v_{j-1}) \approx \frac{1}{k}$. Assim a função $U(v)$ irá utilizar esses cortes para estimar o valor de $W(v)$.

2 O código e a escolha de n e de k

O código escrito em Python tem como objetivo estimar a função $W(v)$ com um erro menor que 0.0005. Desse modo escrevemos algumas funções, que juntas auxiliam a função principal a devolver $U(v)$, o valor da integral calculado até o ponto de corte v .

As funções `x()` e `y()` servem, intuitivamente, para o usuário escolher os vetores x e y e a partir deles a função `theta(n, x, y)` calcula n triplas ordenadas, com distribuição Dirichlet, tal que a soma dos valores dentro de cada tripla é igual a 1. A função `f(theta, x, y)`, usando os vetores x , y e θ , calcula a função de densidade de probabilidade posterior assim como esperado, gerando n valores.

A função `hitu(v1, v2, n, theta, x, y)` usa `f(theta, x, y)` dentro do comando `np.argwhere` para que uma lista seja criada com a posição dos valores que estavam entre v_i e v_{i+1} , posteriormente usando `len` para obter o número total de valores que caíram dentro desse intervalo e dividindo pelo número total de valores n .

A estimativa de $W(v)$ é calculada em si pela função `u()`, que usa todas as funções citadas anteriormente. Após ter chamado `x()` e `y()`, pedimos que o usuário insira o valor v , tal que a integral da função será calculada até ele.

Tendo as entradas e o valor de n e k definido calculamos o vetor θ e aplicamos na função $f(\theta|x, y)$. A função `quantiles` da biblioteca `statistics` gera uma lista com $k - 1$ pontos de corte, pois não inclui os extremos, então adicionamos zero ao começo da lista e o *sup* de f no final, resultando numa lista de comprimento $k + 1$, tal que o número de valores que estão entre v_i e v_{i+1} sobre o número total de pontos é aproximadamente $\frac{1}{k}$. A função `quantiles` divide os valores de $f(\theta|x, y)$ em quantis de igual probabilidade (no nosso caso $\frac{1}{k}$), junto com os extremos, ou seja, a soma de todos os *bins*, deve ser um.

Após termos os cortes definidos, descartamos aqueles que são maiores que v , resultando numa lista de valores de comprimento menor ou igual que a anterior. Como sabemos que o valor da área em cada *bin* é $\frac{1}{k}$, assumimos que a área é $\frac{1}{k}$ vezes o comprimento da nova lista menos 1, ou seja, a soma de todas as parcelas, já que cada parcela tem área igual a $\frac{1}{k}$ e o número de parcelas é de fato o comprimento da nova lista menos 1. Depois checamos se o valor v estava na lista, se não e se o comprimento da lista está menor calculamos a área entre o último valor da lista e v , somando o resultado ao valor da área previamente calculado e assim obtemos a estimativa da integral.

Porém, precisamos de n para realizar esses cálculos. Afim de definir n usamos a função:

$$n = \frac{z^2 \sigma^2}{\epsilon^2}$$

Como o valor da integral vai de zero a um, podemos considerá-la uma proporção, sendo sucesso o número de pontos que caíram dentro do corte e fracasso os que caíram fora. Na função começamos com $n = 100000$ e a partir disso calculamos a estimativa da integral e a variância, usando $p(1 - p)$, sendo p considerado o valor da integral. Tomando $z = 2.575$, para um confiança de 99%, a

variância obtida e tomando o erro como 0.0005, calculamos o novo n . Se o novo n for menor que o n inicial, o mantemos.

Assim que obtivermos o novo n que gera um erro absoluto menor que 0.0005 calculamos a área até o corte v e printamos o resultado.

Também verifica-se necessário a determinação do valor de k . Para esse cálculo percebeu-se que esse método utiliza uma ideia semelhante, senão igual a um histograma de intervalos de largura variável. E para esse tipo de histograma utiliza-se a seguinte relação entre o número de dados n e a quantidade de classes k :

$$k = 4 \left(\frac{2n^2}{z} \right)^{\frac{1}{5}} \quad (1)$$

Assim como na implementação do n , foi colocado um valor inicial para k dentro do código para formar a amostra piloto. E após isso recalcula-se o k a partir da equação (1). Se esse novo k for menor que o k inicial, mantemos o inicial.

3 Resultados e Conclusão

Nessa seção discute-se os resultados obtidos pelo código. Primeiramente, será falado do tempo que a função demorou para rodar, após alguns testes com variando os parâmetros como x , y e k , notou-se um tempo de simulação por volta de 15 segundos, contudo esse tempo varia com os parâmetros e também com o valor de v escolhido. Esse tempo médio foi considerado, pelos integrantes do grupo, um valor aceitável para o tempo. Outro ponto avaliado no código foi o valor do erro da estimativa para a $W(v)$, os valores do erro absoluto, pela fórmula, ficam abaixo de 0.0005, ou seja, atende esse critério. Para verificar a coerência do código foi calculado a $U(v)$ para $v = \sup f(\theta)$, teoricamente o valor encontrado deve ser 1, ao fazer algumas vezes variando os parâmetros, verificou-se que o valor retornado era 1, assim como o esperado, logo pode-se dizer que a função está coerente com sua definição.

Com isso é possível afirmar que a função gera estimativas boas para o valor da função $W(v)$, já que pode-se perceber que o erro absoluto dela fica abaixo de 0.0005. Além disso pode-se dizer que o tempo que leva é, razoavelmente, baixo, logo o código feito segue as condições impostas no enunciado desse EP.