

# EP 4 MAP 2212

Nicholas Gialluca Domene  
N USP 8543417 - IME  
Felipe de Moura Ferreira  
N USP 98674702 - IME  
UNIVERSIDADE DE SÃO PAULO

June 6, 2021

**PROBLEM STATEMENT TAKEN FROM** <https://www.youtube.com/watch?v=ZmaYDTLR8yw>

- Consider the m-dimensional Multinomial statistical model, with observations,  $x$ , prior information,  $y$ , and parameter  $\Theta$ ;

$$- x, y \in N^m, \Theta \in Sm = \{\Theta \in R^{+m} | \Theta'1 = 1\}, m = 3; \quad (1)$$

The statistical model comprises:

- Posterior density potential,

$$f(\theta|x, y) = \frac{1}{B(x+y)} \prod_{i=1}^m \theta_i^{x_i+y_i-1} \quad (2)$$

- Cut-off set,

$$T(v) = \{\Theta \in Sm | f(\theta|x, y) \leq v\}, v \geq 0 \quad (3)$$

- Truth function,

$$W(v) = \int_{T(v)} f(\Theta|x, y) d\Theta \quad (4)$$

$W(v)$  is the posterior probability mass inside  $T(v)$ , that is, the probability mass where the posterior potential,  $f(\Theta|x, y)$ , does not exceed the threshold level  $v$

$$- Dirichlet(\theta|a) = \frac{\prod_{i=1}^m \theta_i^{a_i-1}}{B(a)}, \quad (5)$$

$$B(a) = \frac{\prod_{i=1}^m \Gamma(a_i)}{\Gamma(\sum_{i=1}^m a_i)} \quad (6)$$

$B(a)$  is the multivariate Beta function

- Set  $k$  cut-off points,  $0 = v_0 < v_1 < v_2 < \dots < v_k = \sup f(\Theta)$

- Use a Gamma RNG to generate  $n$  points in  $Sm$ ,  $\Theta_1, \dots, \Theta_n$ , distributed according to the posterior density function

- Use the fraction of simulated points,  $\Theta_t$ , inside each bin;  $v_{j-1} \leq f(\Theta_t) < v_j$ , as an approximation of  $W(v_j) - W(v_{j-1})$

- Dynamically adjust the bin's borders,  $v_j$ , to get bins of approximately equal weight, i.e.,  $W(t_j) - W(t_{j-1}) \approx \frac{1}{k}$

# Implementation

The algorithms were implemented using Python and the usage instructions can be found below.

The overall strategy was treated the implementation as an experiment that using the proposed conditions by the problem statement and a sufficiently large  $n$  (defined below), should yield an estimate to the integral  $W(v)$  with a Confidence Interval lesser or equal to 0.0005 with 95% confidence.

We defined  $n$  through an algebraic manipulation of the Normal approximation of the Binomial distribution:

$$\begin{aligned} Z_{score} &= \frac{mdd}{\sqrt{\frac{\sigma^2}{n}}} \\ Z_{score} \cdot \frac{\sqrt{\sigma^2}}{\sqrt{n}} &= mdd \\ \frac{Z_{score} \cdot \sqrt{\sigma^2}}{mdd} &= \sqrt{n} \\ n &= \frac{Z_{score}^2 \cdot \sigma^2}{mdd^2} \end{aligned} \tag{7}$$

where  $mdd$  is the “Minimal Detectable Difference” which in this case is 0.0005, meaning that in the worst case scenario it will 0.0005 (if  $\int_0^1 f(x)dx = 1$ ). Considering the worst case scenario:

$$mmd = 0.0005, \quad Z_{score} = 1.65, \quad \sigma^2 = 0.25 \tag{8}$$

$$n = \frac{1.65^2 \cdot 0.25}{0.0005^2} = 2722500 \tag{9}$$

## HOW TO RUN:

1. start by initiating the EP4 class defined in `ep4_class.py` inputting the vector  $\mathbf{x}$  and  $\mathbf{y}$ , both of dimension 3. In the `__init__` method, we will store both  $\mathbf{x}$  and  $\mathbf{y}$  vectors, a `alpha` vector that is the sum of each  $x_i$  and  $y_i$  for each  $alpha_i$ , the constant  $B(x+y)$  of the denominator of the evaluating function  $f(\Theta|x,y)$  and the  $n$  for generating  $n$   $\Theta$  observations of the simplex generated by the Dirichlet distribution.

---

```
def __init__(self, x, y):
    self.x = x
    self.y = y
    self.n = 2722500
    self.alpha = [x[i] + y[i] for i in range(len(x))]

    numerator = 1
```

```

for i in range(len(self.alpha)):
    numerator *= math.gamma(self.alpha[i])
B_x_y = numerator / math.gamma(sum(self.alpha))

self.denominator_constant = 1/B_x_y

```

---

2 run the method `generate_theta` that will generate the  $n$   $\Theta$  observations through the Dirichlet distribution with parameters `alpha` and store them into the variable `thetas`.

---

```

def generate_theta(self):
    thetas = np.random.dirichlet(self.alpha, self.n)
    self.thetas = thetas

```

---

3 execute the method `order_f_thetas` that will create a list `f_thetas` with the values of  $f(\Theta, y)$  for each  $\Theta_i$  and sort that list in ascending manner. The method will store the ordered list of  $f(\Theta)$  values, the minimum and the maximum value for  $f(\Theta)$  in the generated observations. This runs in  $O(n \log n)$  time,  $n$  being the number of  $\Theta$  observations generated.

---

```

def order_f_thetas(self):
    f_thetas = [self.f(theta) for theta in self.thetas]
    f_thetas.sort()

    self.ordered_f_thetas = f_thetas
    self.min_f = f_thetas[0] #min value of f_thetas since it is ordered
    self.sup_f = f_thetas[-1] #max value of f_thetas since it is ordered

```

---

4 run the method `U(v)` passing  $v$  as the parameter to the desired output function `U`. The idea behind this method is:

Since we have an ordered list the the values of  $f(\theta)$  for each  $\theta$  in our sample space, then we need to find out how many theta observations have  $f(\theta)$  below certain  $v$ . By finding the index where we would insert a new observation of value  $v$  in our ordered list of  $f(\theta)$ , we use that index to determine how many observations are to the left (lower values). The number of observations whose  $f(\theta)$  values are below  $v$  divided by the total number of observations ( $n$ ) is the estimate for  $W(v)$ . This decreases exponentially the running time and gives a better estimate of  $W(v)$  than using bins since in the worst case scenario, it is exactly the same as using the proposed bins given that the is runs in  $O(\log n)$  time.

---

```

def U(self, v):

    if v > self.sup_f:
        return 1
    if v < self.min_f:
        return 0

    n = self.n

```

```
        i = bisect.bisect_left(self.ordered_f_thetas, v)
        return (i + 1)/n #index divided by total n points
```

---

In the `ep4.py` file, you will find an example of usage passing the `x` and `y` vectors and a few `test_cases`, ready to be tested in the evaluation process:

---

```
from ep4_class import EP4
x = [1, 1, 1]
y = [1, 1, 1]

test_cases = [0, 1, 0.9, 0.00001]

ep4 = EP4(x, y)

ep4.generate_theta()

ep4.order_f_thetas()

for test in test_cases:
    print("U(%s) = %(test), ep4.U(test))
```

---

## Results

---

```
[12:15:38] nicholas.domene ~/Documents/map2212/map2212 (main * u=) $ python3
    ep4.py
Initiating...
U(0) = 0
U(1) = 1
U(0.9) = 1
U(1e-05) = 0.0018938475665748393
Time taken: 8.913041114807129 seconds
```

---