

We want to find the underlying topics of a document, and word-topic associations. Our model:

- Select document length: $N_d \sim \text{Poisson}(\xi)$
- Select topic mixture $\theta_d \sim \text{Dir}(\alpha)$
- For i from 1 to N_d :
 - Select topic $z_{i,d} \sim \text{Multi}(\theta_d)$
 - Select word $w_{i,d} \sim p(w_i|z_{i,d}, \beta_z)$ (i.e. according to topic-specific parameters)

The goal is then to recover the parameters α, β, θ, z . We break this down into two problems:

- 1) *Inference*: Given a new document and known α, β , discover θ and z .
- 2) *Learning*: Given a corpus of documents, estimate appropriate α, β .

We can use estimation methods to achieve this, as follows:

Infer(Document d , α , β):

```

Do:
  # Update word-topic associations
  For  $n$  from 1 to  $N_d$ :
    For  $i$  from 1 to  $K$ :
      # See below for definition of  $\Psi$ 
       $\phi_{n,i}^{(new)} \leftarrow \beta_{i,w_n} \times \exp\left(\Psi\left(\gamma_i^{(old)}\right) - \Psi\left(\sum_j^K \gamma_j^{(old)}\right)\right)$ 
      Normalize( $\phi_n^{(new)}$ )

  # Update topic mixture (here all topics at once)
   $\gamma^{(new)} \leftarrow \alpha + \sum_n \phi_n^{(new)}$ 
Until convergence
Return  $\phi, \gamma$ 

```

Learn(Corpus c):

```

Do:
  # Expectation step
  For document  $d$  in  $c$ :
     $\phi^{(d)}, \gamma^{(d)} \leftarrow \text{Infer}(d, \alpha^{(old)}, \beta^{(old)})$ 
  # Maximization step
   $\beta^{(new)} \leftarrow \sum_d \sum_n \phi_{n,i}^{(d)} w_{d,n}$ 
   $\alpha^{(new)} \leftarrow \text{Newton-Raphson}(\alpha^{(old)})$  #Linear time due to special struct.
Until convergence
Return  $\alpha, \beta$ 

```

#Abramowitz, Stegun. Handbook of Mathematical functions, p. 259

$\Psi(x)$:

```

 $\gamma \leftarrow 0.57721566$  # Euler-Mascheroni constant
 $x' \leftarrow x - 1$ 
sum  $\leftarrow 0$ 
For  $n$  from 1 to many: # Depends on rate of convergence
  sum  $+= x' / (n * (n + x'))$ 
return  $-\gamma + \text{sum}$ 

```