

TEXT MINING FOR SOCIAL SCIENCES

PROBLEM SET 3

Maria Fernández, Tim Kreienkamp, Joan Verdú

2) Consider the following multinomial mixture model for describing documents. First, each document d is allocated a latent topic $z_d \in \{0, \dots, K\}$, where $\Pr[z_d = k] = \rho_k$. Then, each word in document d is generated by taking independent draws from $\beta_{k,j}$, a probability vector with V elements (V is the number of distinct terms in the corpus). Let x_{dv} be the count in document d of term v .

(a) In this model, what are the parameters, what are the latent variables, and what is the observed data?

Parameters are:

ρ_k : Probabilities of a document to belong to a topic (K+1 parameters)

β_k : Vector (dimensions V) of probabilities of distinct terms, for each topic k . At the end there are $(K+1) \cdot V$

Adding all, we have $(K+1)(V+1)$ parameters

Latent variables are the topics of each document (z_d), assuming that each document belongs just to one topic (simple model topic).

Observed data are documents (d , n_d in total) and words (w , V unique ones) of each document

(b) Write down the complete data log-likelihood function. Recall that this expresses the log of the joint probability of the latent variables and the observed data.

The likelihood is a factorization of probabilities of each document. We know the topic of each document is a summation, so the probability of the words of a document is the product of probabilities of each term to the power of counts, weighted by the probability of belonging to the assigned topic (z). That is

$$L(d, w, z; \rho, \beta) = \prod_{d=1}^{n_d} \rho_{d,z_i} \prod_{j=1}^V \beta_{z,j}^{x_{dv}}$$

Then the log likelihood is just (transforming log of product to sums of log):

$$\log L = \sum_{d=1}^{n_d} \log \rho_{d,z_i} \prod_{j=1}^V \beta_{z,j}^{x_{dv}} = \sum_{d=1}^{n_d} \log \rho_{d,z_i} + \sum_{d=1}^{n_d} \sum_{j=1}^V x_{dv} \log \beta_{z,j}$$

Where n_d are the number of documents; ρ_{d,z_i} is the probability of document i of belonging to topic z_i ; then $\beta_{k,j}$ is the probability of term j when topic is z ; and finally recall x_{dv} is the count of term v in document d (so x is function of w and d).

(c) Compute the expected value of the above log-likelihood function given fixed values for the parameters. Recall this is the E-step in the EM algorithm. Denote this function Q , and note that it will depend on the parameter values.

c) E-step

$$Q(\theta_{\text{new}}, \theta') = \sum_{i=1}^{nd} \sum_{z=1}^K P_{\theta'}(z|w, d) \cdot \log P_{\theta}(w, z|d=i)$$

for each doc and each zd (topic) \rightarrow

$$P_{\theta'}(z_d=k|d, w) = \frac{p'_k \cdot P_{\theta'_k}(w_k)}{P_{\theta'}(w_k)} = \frac{p'_k \cdot \prod_{v=1}^V \beta'_{kv}^{x_{iv}}}{\sum_{z=1}^K p'_z \cdot \prod_{v=1}^V \beta'_{kv}^{x_{iv}}}$$

This can be computed with initial $\theta' = p', \beta'$, and knowing x_{iv} (doc-term)

lets call $P_{\theta'}(z|w, d) = Y_{iz}$ (a number depending on doc and value of latent variable (topic))

$$Q(\theta, \theta') = \sum_{i=1}^{nd} \sum_{z=1}^K Y_{iz} \cdot \log p'_{iz} \prod_{v=1}^V \beta'_{kv}^{x_{iv}}$$

\downarrow
unknown known
 $\theta = p, \beta$

(d) Maximize Q with respect to the parameter values. Show ALL steps. (Hint: this will involve the setting up of a constrained maximization problem).

The setup is

$\text{MAX } Q(\theta, \theta')$	\leftrightarrow	$\text{MAX } Q - \lambda_p (\sum_i p_i - 1) - \lambda_{p_1} (\sum_v \beta_{1v} - 1) - \dots - \lambda_{p_K} (\sum_v \beta_{Kv} - 1)$
$\text{st. } \sum_i p_i = 1$		
$\left\{ \sum_v \beta_{iv} = 1 \right\}_{i=1}^K$		

where we also have to consider all $\text{rhos} > 0$ and all $\text{betas} > 0$, since they are probabilities.

Maximizing the rhos we have:

d) M-step scalar, defined at E-step

optimize all p, β

$$(1) \quad \frac{\partial Q}{\partial p_z} = \sum_{i=1}^{nd} Y_{iz} \cdot \frac{\prod_{v=1}^V \beta_{zv}^{x_{iv}}}{p_z \cdot \prod_{v=1}^V \beta_{zv}^{x_{iv}}} - \lambda_p = -\lambda_p + \frac{1}{p_z} \sum_{i=1}^{nd} Y_{iz} = 0$$

$\} \text{ for each } z=1 \dots K \text{ (topics)}$

$\ln u = \frac{w}{u}$

$-R \text{ (a scalar)}$

Max betas:

$$(2) \quad \frac{\partial Q}{\partial \beta_{zt}} = \sum_{i=1}^{nd} Y_{iz} \cdot \frac{X_{it} \cdot \beta_{zt} \cdot P_z \cdot \prod_{v=1}^V \beta_{zv}}{\prod_{v=1}^V \beta_{zv}} - \lambda_{\beta_z} =$$

$\ln u = \frac{w}{u}$

$= R_{zt}$ (a scalar)

for each β_{kv}
Topic: $v \in \{1, \dots, V\}$
Term: $t \in \{1, \dots, T\}$

$$= \sum_{i=1}^{nd} Y_{iz} \frac{X_{it} \cdot \beta_{zt} \cdot R_{zt}}{\beta_{zt} \cdot R_{zt}} - \lambda_{\beta_z} =$$

$$= -\lambda_{\beta_z} + \sum_{t=1}^T \frac{Y_{it} \cdot X_{it}}{\beta_{zt}} = -\lambda_{\beta_z} + \frac{1}{\beta_{zt}} \sum_{i=1}^{nd} Y_{iz} \cdot X_{it} = 0$$

And max lambdas:

$$(3) \quad \frac{\partial Q}{\partial \lambda_i} = 0 \rightarrow \text{recover constraints}$$

$$\sum_z P_z = 1$$

$$\text{For each topic: } \sum_v \beta_{zv} = 1$$

So we have a solvable system of linear equations.

(SUMMARY : (1) K equations <u>LINEAR EQUATIONS</u> (2) $V \cdot K$ equations (3) $1+K$ equations	<u>UNKNOWNs:</u> $P_z: K$ $\beta_{zv}: V \cdot K$ $\lambda: 1+K$
<hr/>	<hr/>
$(V+2) \cdot K + 1$ eqs	$(V+2) \cdot K + 1$ unknown

THIS IS A SYSTEM OF LINEAR EQUATIONS WITH SAME # OF UNKNOWN \rightarrow SOLVABLE
(SIMPLEX FOR INSTANCE)

(e) Use your answers to the above to write pseudo-code for implementing the EM algorithm for multinomial mixture models.

Basically to find a local optimum for the parameters rho and beta we have to do iteratively steps E and M of the EM algorithm until the variation of log L gets under a predefined threshold.

```

# Load/calculate data
Docs= load_docs()
Words= load_words()
Topics=load_topics()
Z= unique(topics)
D= unique(Docs)
W= unique(words)
K=length(Z)
Nd=length(D)
V=length(W)
X = count(docs,words) #matrix of Nd x V
# Initialize parameters and
Rho= random(K)
Beta= matrix(random(K*V),rows=K,cols=V)
Diff= 1
Delta= 0.001
P= matrix(NA,rows=Nd,cols=K)
# Iterate EM algorithm
Max_Iter=1000
Iter=0
While (abs(diff)>delta AND iter<Max_Iter) do
    #E step
    S=0
    For (d in Docs) :
        Denominator= for (z in Z) :
            S=S+rho(z)*factorial(beta(row=z),Xdv)
        For (z in Z) :
            P[d,z]=[rho(z)*factorial(beta(row=z),Xdv)]/ denominator
    # we only need this to go for optimization, the other part of Q is fixed
    #M step
    Rho_old=rho
    Beta_old=beta
    #solve linear equations to update rho, beta values
    Rho,Beta= solve (P,Rho_old,Beta_old,Xdv,)
    #Compute gain
    Diff=1-[LogL(rho,beta,xdv)/LogL(rho_old,beta_old_xdv)]
    Iter += 1

```

Since this is dependent on initial values of parameters, to be closer to the global optimum, we should repeat the procedure several times, and get the optimal solution than yields maximum value of log L.

If we consider previous code as function calculate_EM that yields estimated rho, beta coefficients and optimum logL, we have this pseudocode (considering data includes docs,words and topics):

```
# Repeat EM and get best result
Nsampling=100
RhoF,BetaF,LogLF=calculate_EM(data)
For (I in 1:Nsampling) :
    Rho,beta,logL=calculate_EM(data)
    If logL>logLF :
        rhoF=rho
        betaF=beta
        logLF=logL
```