

# Non-convex optimization in RNA-seq models

*Kushal K Dey*

*December 29, 2015*

## Overview

RNA-seq is a popularly used technique in genetics/genomics to map reads (small nucleotide sequences) coming from a sample (may be single cell or a tissue sample) to the genome. The number of reads mapped to each gene is recorded, resulting in counts data. The number of samples or repetitions of this experiment are in the range of 100 to 1000, whereas the number of genes range from 20,000 to 50,000. The reads data may be statistically analyzed as it is (Poisson framework for counts) or may be normalized to RPKM values which is amenable to continuous distribution fitting (typically Normal framework).

Under topic model settings, we use a mixed membership approach to determine the proportion of reads coming from different cell subpopulations (tissue level data) or cell phases (single cell data). In simple words, we wish to soft cluster the samples into groups and then record the proportional membership of the sample to each cluster or group.

## Model

**Reads model** If we model the counts of reads, then we may use the following Poisson framework. Let  $c_{ng}$  be the number of counts of reads from sample  $n$  that get mapped to gene  $g$ .

MODEL 1 - Identity link model

$$c_{ng} \sim \text{Poi}\left(\sum_{k=1}^K \omega_{nk} \exp(\alpha_{kg})\right) \quad \sum_{k=1}^K \omega_{nk} = 1$$

MODEL 2 - log link model

$$c_{ng} \sim \text{Poi}\left(\exp\left(\sum_{k=1}^K \omega_{nk} \alpha_{kg}\right)\right) \quad \sum_{k=1}^K \omega_{nk} = 1$$

We may assume the priors for  $\omega$  as follows

CHOICE 1

$$(\omega_{n1}, \omega_{n2}, \dots, \omega_{nK}) \sim \text{Dir}_K\left(\frac{1}{K}, \frac{1}{K}, \dots, \frac{1}{K}\right)$$

CHOICE 2

$$(\omega_{n1}, \omega_{n2}, \dots, \omega_{nK}) \sim \text{Dir}_K(1, 1, \dots, 1)$$

We assume the prior for  $\alpha$  to be as follows

$$\alpha_{kg} \sim N(0, 1000)$$

This is not fully non-informative theoretically but close to being non-informative practically and easier to handle as well.

The MODEL 1 or the identity link model is biologically more meaningful whereas MODEL 2 is basically the topic model equivalent of the generalized linear model under Poisson family with log link, and easier to estimate using stepwise iterative estimation of parameters.

**RPKM model** If we are provided with the RPKM values, which are log transformed and scaled values obtained from the reads or counts values, we can fit the following model.

Let  $r_{ng}$  be the RPKM for sample  $n$  and gene  $g$ . Then we can write

$$r_{ng} \sim N\left(\sum_{k=1}^K \omega_{nk} \alpha_{kg}\right)$$

with the prior choices as reported above.

## Simulation experiments

I am providing the simulated data for the three models so that they are ready to be applied the model and tested for whether the method works on these examples (the estimated  $\omega$  and  $\alpha$  should match with the true values from which the data has been simulated from).

### MODEL 1 - Identity link model

```
K=4;
G=100;
N=500;

alpha_true=matrix(rnorm((K)*G,0.5),nrow=(K)); ### the matrix of fixed effects

library(gtools)
T=10;
omega_true=matrix(rbind(rdirichlet(T*10,c(3,4,2,6)),rdirichlet(T*10,c(1,4,6,3)),
                        rdirichlet(T*10,c(4,1,2,2)),rdirichlet(T*10,c(2,6,3,2)),
                        rdirichlet(T*10,c(3,3,5,4))), nrow=N);

### generating the table

read_counts_model1=matrix(0,N,G);
for(n in 1:N)
{
  for(g in 1:G)
  {
    mean=omega_true[n,]*%exp(alpha_true[,g]);
    read_counts_model1[n,g]=rpois(1,mean);
  }
}
```

## MODEL 2 - Log link model

```
K=4;
G=100;
N=500;

alpha_true=matrix(rnorm((K)*G,0.5),nrow=(K)); ### the matrix of fixed effects

library(gtools)
T=10;
omega_true=matrix(rbind(rdirichlet(T*10,c(3,4,2,6)),rdirichlet(T*10,c(1,4,6,3)),
                        rdirichlet(T*10,c(4,1,2,2)),rdirichlet(T*10,c(2,6,3,2)),
                        rdirichlet(T*10,c(3,3,5,4))), nrow=N);

### generating the table

read_counts_model2=matrix(0,N,G);
for(n in 1:N)
{
  for(g in 1:G)
  {
    mean=exp(omega_true[n,]*%*%alpha_true[,g]);
    read_counts_model2[n,g]=rpois(1,mean);
  }
}
```

## MODEL 3 - normal RPKM model

```
K=4;
G=100;
N=500;

alpha_true=matrix(rnorm((K)*G,0.5),nrow=(K)); ### the matrix of fixed effects

library(gtools)
T=10;
omega_true=matrix(rbind(rdirichlet(T*10,c(3,4,2,6)),rdirichlet(T*10,c(1,4,6,3)),
                        rdirichlet(T*10,c(4,1,2,2)),rdirichlet(T*10,c(2,6,3,2)),
                        rdirichlet(T*10,c(3,3,5,4))), nrow=N);

### generating the table

rpkm_model3=matrix(0,N,G);
for(n in 1:N)
{
  for(g in 1:G)
  {
    mean=omega_true[n,]*%*%alpha_true[,g];
    rpkm_model3[n,g]=rnorm(1,mean);
  }
}
```

```
}  
}
```

## Discussion

- Of the parameters, only  $\omega$  are constrained, they sum to 1 for each  $n$ . The idea would be to de-simplex them using

$$\nu_{nk} = \log\left(\frac{\omega_{nk} + 1e-06}{\omega_{n1} + 1e-06}\right) \quad k = 2, 3, \dots, K$$

Carry out the optimization on  $\nu_{nk}$  and then reverse transform them.

- The time and computational expense is an important issue in this model fit. Note that the real data model will have 100 s of samples and 20,000 to 50,000 genes. So, need to see how fast the algorithm is on this toy dataset provided above, and how it scales for increase in number of samples and genes. So, report the time expenditure (using *system.time()*). We may need to use variable selection methods to select genes if it is found to be slow.
- We shall apply the method on real data (both reads and RPKM) once the method has been validated on the above simulation examples.
- The model may have overdispersion effects and also possible technical effects, which would make the model slightly more complicated, but we are not worrying about these as of now. First, we focus on making the method work on the simple set ups.