

# Phenotyping via Bayesian Nonparametric Tensor Factorization

**Abstract—We do phenotyping.**

## I. INTRODUCTION

## II. MODEL

### A. CANDECOMP/PARAFAC (CP) Tensor Decomposition

A tensor  $T \in R^{n_1 \times n_2 \times \dots \times n_d}$  can be decomposed by a linear combination of rank-1 tensors:

$$T = \sum_{i=1}^r \lambda_i a_i^1 \otimes a_i^2 \dots a_i^d$$

Where  $\lambda_i \in R$ ,  $a_i^k \in R^{n_k}$  and  $\otimes$  represents the outer product. The  $r$  is the rank of tensor when  $r$  is minimized in above equation. When  $r$  is nor minimal, then the above decomposition is referred to as CANDECOMP/PARAFAC decomposition.

### B. Dirichlet Process Tensor Factorization

TABLE I. DATA PARAMETERS

Notation	Meaning
$\lambda$	The weight vector of phenotypes, here $\ \lambda\  = 1$
$\lambda_k$	The weight of $k$ -th phenotype; $\lambda_k \in [0, 1]$
$N_p$	Number of patients
$N_m$	Number of medications
$N_d$	Number of diagnosis
$T$	Record Tensor. $T \in \{0, 1\}^{N_p \times N_m \times N_d}$

TABLE II. MODEL PARAMETERS

Notation	Meaning
$\alpha$	The hyper parameter of Beta Distribution
$\beta_i$	The stick-breaking process random variable. $\beta_i \sim Beta(1, \alpha)$
$\gamma_p$	The hyper parameter of prior Dirichlet Distribution of patient
$\gamma_m$	The hyper parameter of prior Dirichlet Distribution of medication
$\gamma_d$	The hyper parameter of prior Dirichlet Distribution of diagnosis
$\theta^{(p,k)}$	The multinomial distribution over patients given the phenotype $k$
$\theta^{(m,k)}$	The multinomial distribution over medications given the phenotype $k$
$\theta^{(d,k)}$	The multinomial distribution over diagnosis given the phenotype $k$

We place a stick-breaking process prior on the  $\lambda$ . The parameter  $\beta_i$  follows the Beta Distribution  $\beta_i \sim Beta(1, \alpha)$ ,  $i = 1, 2, \dots$ . And:

$$\lambda_k = \beta_k \prod_{i=1}^{k-1} (1 - \beta_i), k = 1, 2, \dots$$

A  $k$ -dimensional Dirichlet random variable  $\theta$  can take values in the  $(k-1)$ -simplex where  $\sum_{i=1}^k \theta_i = 1$ , and the probability density is

$$p(\theta|\gamma) = \frac{\mathcal{T}(\sum_{i=1}^k \gamma_i)}{\prod_{i=1}^k \mathcal{T}(\gamma_i)} \theta_1^{\gamma_1-1} \dots \theta_k^{\gamma_k-1}$$

It is natural to place the conjugate prior (Dirichlet Distribution) over the multinomial distributions parameterized by  $\theta^{(p)}$ ,  $\theta^{(m)}$  and  $\theta^{(d)}$ :

$$\theta^{(p,k)} \sim Dir(\gamma_p), k = 1, 2, \dots$$

$$\theta^{(m,k)} \sim Dir(\gamma_m), k = 1, 2, \dots$$

$$\theta^{(d,k)} \sim Dir(\gamma_d), k = 1, 2, \dots$$

The patients, medications and diagnosis all follow the multinomial distribution parameterised by  $\theta$ , specifically:

$$\text{patient}_i \sim Multi(\theta^{(p)}), i = 1, 2, \dots, N_p$$

$$\text{medication}_i \sim Multi(\theta^{(m)}), i = 1, 2, \dots, N_m$$

$$\text{diagnosis}_i \sim Multi(\theta^{(d)}), i = 1, 2, \dots, N_d$$

Given the parameters, the likelihood of given record triplet  $(i, j, k)$  (patient  $i$  had diagnostic record  $k$ , and has taken medicine  $j$ ) can be formulated as:

$$p(T_{i,j,k}|\alpha, \gamma_p, \gamma_m, \gamma_d) = \sum_{r=1}^{\infty} \lambda_r p(i|r, \gamma_p) p(j|r, \gamma_m) p(k|r, \gamma_d)$$

## III. EXPERIMENTS

## IV. CONCLUSION

## ACKNOWLEDGMENT

## REFERENCES

- [1] H. Kopka and P. W. Daly, *A Guide to L<sup>A</sup>T<sub>E</sub>X*, 3rd ed. Harlow, England: Addison-Wesley, 1999.