

# TENSOR NORMAL MODELS ON DIRECTED ACYCLIC GRAPHS

ABSTRACT. Abstract here.

## 1. INTRODUCTION

Maximum Likelihood Estimation (MLE) is a fundamental problem in statistics, and recently it has been studied from the lens of algebra. The MLE problem is to find a point in a model that ‘best’ fits some data. The idea of ‘fitting’ a data is given by a *likelihood function*, and the problem of finding a point that ‘best’ fits the data translates to *maximizing* the likelihood function. A point that maximizes the likelihood function is called the *maximum likelihood estimate*.

Throughout we will use the calligraphic font  $\mathcal{G}$  to denote a directed acyclic graph (DAG) and the normal font  $G$  to denote a group.  $\mathcal{M}$ , with some subscript containing an object, will be used to denote a model. For example,  $\mathcal{M}_G$  is the model associated to a group  $G$ .  $\text{PD}_m$  will refer to the cone of symmetric  $m \times m$  positive definite matrices.

## 2. MAIN RESULTS

The main paper [1] deals mainly about Gaussian group models and section 5 talks about certain DAG models and their relationship with Gaussian group models. Anna Seigal suggested to deal with the model  $\mathcal{M}(\mathcal{G}_1, \mathcal{G}_2) := \{\Psi_1 \otimes \Psi_2 : \Psi_i \in \mathcal{M}_{\mathcal{G}_i}\}$  when two DAGs  $\mathcal{G}_1, \mathcal{G}_2$  are given, and raised the primary question that when are such models Gaussian group models, that is, when is  $\mathcal{M}(\mathcal{G}_1, \mathcal{G}_2) = \mathcal{M}_G$  for some group  $G$  of matrices. For the one-parameter case, such models have been studied. Given a DAG  $\mathcal{G}$ , one considers the set of matrices

$$G(\mathcal{G}) = \{X \in GL_{|V(\mathcal{G})|} : X_{ij} = 0 \text{ for } i \neq j \text{ with } j \not\rightarrow i \text{ in } \mathcal{G}\}.$$

This is relevant in the ‘good’ case when  $\mathcal{M}_{\mathcal{G}} = \mathcal{M}_{G(\mathcal{G})}$  as indicated in theorem 3.2.

The problem on  $\mathcal{M}(\mathcal{G}_1, \mathcal{G}_2)$  naturally starts by considering

$$G(\mathcal{G}_1, \mathcal{G}_2) := \{\Psi_1 \otimes \Psi_2 : \Psi_i \in G(\mathcal{G}_i)\} = G(\mathcal{G}_1) \otimes G(\mathcal{G}_2).$$

It can be directly proven that

**Proposition 2.1.**  $G(\mathcal{G}_1, \mathcal{G}_2)$  is a group iff both  $\mathcal{G}_i$  are TDAGs. If both  $\mathcal{G}_i$ 's are TDAGs, the model  $M = \{\Psi_1 \otimes \Psi_2 : \Psi_i \in \mathcal{M}_{\mathcal{G}_i}\}$  is exactly  $\mathcal{M}_{G(\mathcal{G}_1, \mathcal{G}_2)}$ .

My main contribution is to define a construction  $\mathcal{G}_1 \otimes \mathcal{G}_2$  such that  $G(\mathcal{G}_1) \otimes G(\mathcal{G}_2) \cong G(\mathcal{G}_1 \otimes \mathcal{G}_2)$ . In fact such a construction is commutative (upto relabelling) and associative. Thus it extends to  $\mathcal{G}_1 \otimes \cdots \otimes \mathcal{G}_n$  so that  $\bigotimes_{i=1}^n G(\mathcal{G}_i) \cong G\left(\bigotimes_{i=1}^n \mathcal{G}_i\right)$ . That is, this construction extends to the tensor normal models on  $n$  DAGs.

### 3. BACKGROUND

**3.1. Maximum likelihood estimation.** We will focus on multivariate Gaussian distributions of mean zero and covariance matrix  $\Sigma$ , whose density is given by

$$f_{\Sigma}(\mathbf{y}) = \frac{1}{\sqrt{\det(2\pi\Sigma)}} \exp\left(-\frac{1}{2}\mathbf{y}^{\top}\Sigma^{-1}\mathbf{y}\right)$$

where  $\mathbf{y} \in \mathbb{R}^m$  and  $\Sigma \in \text{PD}_m$ . The corresponding  $\Psi = \Sigma^{-1}$  will be its *concentration matrix*. A *Gaussian model* is a set of  $m$ -dimensional Gaussian distributions with mean zero. Such a model is given by a set of  $m \times m$  symmetric positive definite covariance matrices. Equivalently they are also determined by a set of concentration matrices in  $\text{PD}_m$ . We will refer to this as the Gaussian model, instead of the set of densities themselves. So a Gaussian model  $\mathcal{M}$  is just a subset of  $\text{PD}_m$  whose elements are to be thought of as concentration matrices.

A maximum likelihood estimate is a point  $\Psi$  in the model which maximizes the likelihood of observing some sample data point  $\vec{\mathbf{Y}} = (\mathbf{y}_1, \dots, \mathbf{y}_n)$ . Mathematically, we want a  $\Psi \in \mathcal{M}$  that maximizes the likelihood function  $L_{\vec{\mathbf{Y}}}(\Psi) = \prod_{i=1}^n f_{\Psi^{-1}}(\mathbf{y}_i)$ . Often, it is easier to maximize  $l_{\vec{\mathbf{Y}}} := \log L_{\vec{\mathbf{Y}}}$

instead of  $L_{\tilde{\mathbf{Y}}}$  itself, and they have the same maximizers. Note that

$$\begin{aligned}
 l_{\tilde{\mathbf{Y}}}(\Psi) &= \log L_{\tilde{\mathbf{Y}}}(\Psi) \\
 &= \sum_{i=1}^n \left( -\frac{1}{2} \log \det(2\pi\Psi^{-1}) - \frac{1}{2} \mathbf{y}_i^\top \Psi \mathbf{y}_i \right) \\
 \implies \frac{2}{n} l_{\tilde{\mathbf{Y}}} &= -\log \det(2\pi\Psi^{-1}) - \frac{1}{n} \sum_{i=1}^n \mathbf{y}_i^\top \Psi \mathbf{y}_i \\
 &= -m \log(2\pi) + \log \det \Psi - \frac{1}{n} \sum_{i=1}^n \mathbf{y}_i^\top \Psi \mathbf{y}_i \\
 \implies \frac{2}{n} l_{\tilde{\mathbf{Y}}} + m \log(2\pi) &= \log \det \Psi - \frac{1}{n} \sum_{i=1}^n \text{Tr} \left( \mathbf{y}_i^\top \Psi \mathbf{y}_i \right) \\
 &= \log \det \Psi - \text{Tr} \left( \frac{1}{n} \sum_{i=1}^n \Psi \mathbf{y}_i^\top \mathbf{y}_i \right) \\
 &= \log \det \Psi - \text{Tr} \left( \Psi \frac{1}{n} \sum_{i=1}^n \mathbf{y}_i \mathbf{y}_i^\top \right) \\
 &= \log \det \Psi - \text{Tr} (\Psi S_{\tilde{\mathbf{Y}}})
 \end{aligned}$$

where  $S_{\tilde{\mathbf{Y}}} = \frac{1}{n} \sum_{i=1}^n \mathbf{y}_i \mathbf{y}_i^\top$  is the sample covariance matrix. This is clearly symmetric positive semi-definite. The above calculation shows that  $\arg\max_{\Psi \in \mathcal{M}} \{l_{\tilde{\mathbf{Y}}}(\Psi)\} = \arg\max_{\Psi \in \mathcal{M}} \{\log \det \Psi - \text{Tr}(\Psi S_{\tilde{\mathbf{Y}}})\}$ . Take  $\ell_{\tilde{\mathbf{Y}}}(\Psi) := \log \det \Psi - \text{Tr}(\Psi S_{\tilde{\mathbf{Y}}})$ . Observe that  $\ell_{\tilde{\mathbf{Y}}}(\Psi)$  is invariant under similarity of  $\Psi$  and similarity of  $S_{\tilde{\mathbf{Y}}}$ . Further  $\Psi, S_{\tilde{\mathbf{Y}}}$  are real symmetric matrices, hence have real eigenvalues and are diagonalizable. This means that if  $\{\lambda_i\}_{i=1}^m \subseteq \mathbb{R}^{>0}$  and  $\{\mu_i\}_{i=1}^m \subseteq \mathbb{R}^{\geq 0}$  are eigenvalues of  $\Psi$  and  $S_{\tilde{\mathbf{Y}}}$  respectively, then  $\ell_{\tilde{\mathbf{Y}}}(\Psi) = \sum \log \lambda_i - \sum \lambda_i \mu_i$ . Note that if  $S_{\tilde{\mathbf{Y}}}$  is invertible,

then

$$\begin{aligned}
\ell_{\vec{Y}}(\Psi) &= \log \det \Psi - \text{Tr}(\Psi S_{\vec{Y}}) \\
&= \log \det \Psi - \log \det S_{\vec{Y}}^{-1} + \log \det S_{\vec{Y}}^{-1} \\
&\quad - \text{Tr}(\Psi S_{\vec{Y}}) + \text{Tr} \mathbf{1}_m - \text{Tr}(S_{\vec{Y}}^{-1} S_{\vec{Y}}) \\
&= \sum \log(\lambda_i \mu_i) + \log \det S_{\vec{Y}}^{-1} - \sum \lambda_i \mu_i + m - \text{Tr}(S_{\vec{Y}}^{-1} S_{\vec{Y}}) \\
&= \sum [\log(\lambda_i \mu_i) - \lambda_i \mu_i] + m + \ell_{\vec{Y}}(S_{\vec{Y}}^{-1}) \\
&\leq \sum_{i=1}^m (-1) + m + \ell_{\vec{Y}}(S_{\vec{Y}}^{-1}) = \ell_{\vec{Y}}(S_{\vec{Y}}^{-1})
\end{aligned}$$

If  $S_{\vec{Y}}$  is singular, assume WLOG  $\mu_1 = 0$ , then  $\ell_{\vec{Y}}(\Psi) = \log \lambda_1 + \sum_{i \geq 2} (\log \lambda_i - \lambda_i \mu_i)$  diverges to  $+\infty$  as  $\lambda_1 \rightarrow \infty$ .

**3.2. Gaussian group models.** The *Gaussian group model* given by a group  $G \subseteq \text{GL}(\mathbb{R}^n)$  is the multivariate Gaussian model comprising all densities of mean zero and concentration matrices given by

$$\mathcal{M}_G := \left\{ X^\top X : X \in G \right\}.$$

The importance of such models is that finding an MLE is equivalent to an optimization problem. This is seen by the following calculation

$$\begin{aligned}
(1) \quad -\ell_{\vec{Y}}(X^\top X) &= \frac{1}{n} \sum \text{Tr}((X \mathbf{y}_i)^\top X \mathbf{y}_i) - \log \det(X^\top X) \\
(2) \quad &= \frac{1}{n} \left\| X \cdot \vec{Y} \right\|_2^2 - \log(\det X)^2
\end{aligned}$$

So maximizing  $\ell_{\vec{Y}}(X^\top X)$  is equivalent to minimizing  $\frac{1}{n} \left\| X \vec{Y} \right\|_2^2 - \log \det(X^\top X)$  which consists of minimizing norms. The exact statement is captured in the following theorem in [1, Proposition 3.4]:

**Proposition 3.1.** *Let  $\vec{Y} \in V^n$  be a tuple of samples. If the group  $G \subseteq \text{GL}(V)$  is closed under non-zero scalar multiples, the supremum of the log-likelihood  $\ell_{\vec{Y}}(\Psi)$  over  $\mathcal{M}_G$  is the double infimum*

$$-\inf_{\lambda > 0} \left( \frac{\lambda}{n} \left( \inf_{H \in G_{\text{SL}}^\pm} \left\| H \cdot \vec{Y} \right\|_2^2 \right) - \dim(V) \log \lambda \right).$$

The MLEs, if they exist, are the matrices  $\frac{n \dim V}{\|H \cdot \vec{Y}\|^2} H^\top H$ , where  $H$  minimizes  $\|H \cdot \vec{Y}\|$  under the action of  $G_{\text{SL}}^\pm$  on  $V^n$ .

In the above,  $G_{\text{SL}}^\pm = \{X \in G : \det X = \pm 1\}$ .

What the above proposition says is that  $\ell_{\vec{Y}}(X^\top X)$  can be maximized in two steps:

- (1) Minimize the norm  $\|H \cdot \vec{Y}\|^2$  over  $H \in G_{\text{SL}}^\pm$ .
- (2) Find a scalar  $\mu \in \mathbb{R}$  so that  $X := \mu H$  minimizes  $-\ell_{\vec{Y}}(X^\top X)$  in eq. (2).

*Proof.* We want to minimize the function  $f : G \rightarrow \mathbb{R}$  given by

$$f(X) = \frac{1}{n} \|X \cdot \vec{Y}\|_2^2 - \log(\det X)^2.$$

Let  $m = \dim V$ . Observe that  $f|_{G_{\text{SL}}^\pm}$  determines  $f$  completely. This is because for any  $X \in G$ , take  $Z := \frac{1}{\mu_X} X$ , where  $\mu_X = (|\det X|)^{\frac{1}{m}} \in \mathbb{R}_{>0}$ . Clearly  $\det Z = \frac{\det X}{|\det X|} = \pm 1$ . Then

$$\begin{aligned} f(X) &= \frac{\mu_X^2}{n} \|Z \cdot \vec{Y}\|_2^2 - \log((\mu_X^m \det Z)^2) \\ &= \frac{\mu_X^2}{n} \|Z \cdot \vec{Y}\|_2^2 - 2m \log \mu_X \\ &= \mu_X^2 f(Z) - 2m \log \mu_X \end{aligned}$$

Note that for  $K > 0$ , the function  $s \mapsto sK - \log s$  minimizes at  $s = \frac{1}{K}$  giving a min value of  $1 + \log K$ , the latter being an increasing function of  $K$ . Thus minimizing  $f(X)$  is equivalent to first minimizing the norm in the orbit of  $G_{\text{SL}}^\pm$  and then minimizing the overall expression with the previous minima. In other words,

$$\inf_{X \in G} f(X) = \inf_{\mu > 0} \left( \mu^2 \cdot \inf_{Z \in G_{\text{SL}}^\pm} \|Z \cdot \vec{Y}\|_2^2 - m \log \mu^2 \right).$$

Replacing  $\lambda = \mu^2$  gives the expression in the expression in the proposition. The MLE, if it exists, is given by the point  $\hat{X} = \hat{\mu} \hat{Z}$ , where

$\hat{Z}$  minimizes  $\left\|Z \cdot \vec{Y}\right\|_2^2$  and  $\hat{\mu} = \frac{\sqrt{mn}}{\left\|\hat{Z} \cdot \vec{Y}\right\|_2}$ , which corresponds to the input  $\hat{\Psi} = \hat{X}^\top \hat{X} = \hat{\mu}^2 \hat{Z}^\top \hat{Z}$ .  $\blacksquare$

**3.3. Transitive DAGs.** The relevance of transitive DAGs in the statistical context was introduced in [1, §5]. We briefly introduce the relevant details here. Every DAG  $\mathcal{G}$  has a model associated to them given by

$$\mathcal{M}_{\mathcal{G}} := \left\{ (\mathbf{1} - \Lambda)^\top \Omega^{-1} (\mathbf{1} - \Lambda) : \Lambda, \Omega \in \mathbb{R}^{m_1 \times m_i}, \Lambda_{ij} \neq 0 \implies j \rightarrow i, \Omega \text{ diagonal and PD} \right\}.$$

Note that taking  $\Omega = \mathbf{1}, \Lambda = \mathbf{0}$  forces  $\mathbf{1} \in \mathcal{M}_{\mathcal{G}}$ . These matrices  $\Psi \in \mathcal{M}_{\mathcal{G}}$  are to be thought of as concentration matrices, so they define a model  $\{f_{\Psi^{-1}}\}$  (which is just a collection of probability densities) given by

$$f_{\Sigma}(y) = \frac{1}{\sqrt{\det(2\pi\Sigma)}} \exp\left(-\frac{1}{2}y^\top \Sigma^{-1}y\right)$$

where  $\Sigma = \Psi^{-1}$ .

$\mathcal{G}$  is said to be a *transitive* DAG (TDAG) if  $i \rightarrow k$  is an edge in  $\mathcal{G}$  whenever  $i \rightarrow j, j \rightarrow k$  are. They turn out to be the ‘good’ DAGs that help relate these models to Gaussian group models due to the following proposition in [1]:

**Theorem 3.2.**  *$G(\mathcal{G})$  is a group iff  $\mathcal{G}$  is a TDAG. In such a case,  $\mathcal{M}_{\mathcal{G}} = \mathcal{M}_{G(\mathcal{G})}$ .*

3.4. ss.

## REFERENCES

- [1] Carlos Améndola, Kathlén Kohn, Philipp Reichenbach, and Anna Seigal. Invariant theory and scaling algorithms for maximum likelihood estimation. *SIAM Journal on Applied Algebra and Geometry*, 5(2):304–337, Jan 2021.