

# Alternating Least Squares in Generalized Linear Models

Arshak Minasyan

Yerevan State University

`arsh.minasyan@gmail.com`

February 23, 2022

## Abstract

We derived a convergence result for a sequential procedure known as alternating maximization (minimization) to the maximum likelihood estimator for a pretty large family of models - Generalized Linear Models (GLMs). Alternating procedure for linear regression becomes to the well-known algorithm of Alternating Least Squares (ALS), because of the quadraticity of log-likelihood function  $L(\boldsymbol{v})$ . In GLMs framework we lose quadraticity of  $L(\boldsymbol{v})$ , but still have concavity due to the fact that error-distribution is from exponential family (EF). Concentration property makes the Taylor approximation of  $L(\boldsymbol{v})$  up to the second order accurate and makes possible the use of alternating minimization (maximization) technique. Examples and experiments confirm convergence result followed by the discussion of the importance of initial guess.

# 1 Introduction

Many statistical tasks can be viewed as problems of semi-parametric estimation when the unknown data distribution is described by a high or infinite dimensional parameter while the target is of low dimension. Typical examples are provided by functional estimation, estimation of a function at a point, or simply by estimating a given sub-vector of the parameter vector. The classical statistical theory provides a general solution to this problem: estimate the full parameter vector by the maximum likelihood method and project the obtained estimate onto the target subspace. This approach is known as *profile maximum likelihood* and it appears to be *semi-parametrically efficient* under some mild regularity conditions, which in case of Generalized Linear Models are satisfied. For more general case, for example, in M-estimation framework these technical conditions should be introduced separately and checked whether they are fulfilled or not. We refer to [8], [6] and the book of Kosorok [10] for a detailed presentation.

This study revisits the problem of profile semi-parametric estimation (see [4], [5] and references therein). One issue that is worth mentioning is the model mis-specification. In most of the cases of practical problems, it is unrealistic to expect that the model assumptions are exactly fulfilled, even if some rich non-parametric family is used. This means that the true data distribution  $\mathbb{P}$  does not belong to the considered parametric family, in our case — exponential family. Applicability of the general semi-parametric theory in such cases is questionable. An important feature of the presented approach is that it equally applies under a possible model mis-specification.

Consider the following statistical model that assumes the unknown data distribution  $\mathbb{P}$  belongs to a given parametric family  $(\mathbb{P}_{\mathbf{v}})$ :

$$\mathbf{y} \sim \mathbb{P} = \mathbb{P}_{\mathbf{v}^*} \in (\mathbb{P}_{\mathbf{v}}, \mathbf{v} \in \Theta), \quad (1.1)$$

where  $\Theta$  is some high dimensional or even infinite dimensional parametric space.

The maximum likelihood approach in the parametric estimation suggests to estimate the whole parameter vector  $\mathbf{v}$  by maximizing the corresponding log-likelihood

$$L(\mathbf{v}) = \log \frac{d\mathbb{P}_{\mathbf{v}}}{d\mu_0}$$

for some dominating measure  $\mu_0$ . Define the maximum likelihood estimator in the following way

$$\tilde{\mathbf{v}} \stackrel{\text{def}}{=} \arg \max_{\mathbf{v} \in \Theta} L(\mathbf{v}). \quad (1.2)$$

Our study admits a model specification  $\mathbb{P} \notin (\mathbb{P}_{\mathbf{v}}, \mathbf{v} \in \Theta)$ . Equivalently, one can say that  $L(\mathbf{v})$  is the *quasi log-likelihood function* on  $\Theta$ . The *target* value  $\mathbf{v}^*$  of the parameter  $\mathbf{v}$  can be defined by

$$\mathbf{v}^* \stackrel{\text{def}}{=} \arg \max_{\mathbf{v} \in \Theta} \mathbb{E} L(\mathbf{v}). \quad (1.3)$$

Under model misspecification,  $\mathbf{v}^*$  defines the best parametric fit to  $\mathbb{P}$  by the considered family. For the results of similar kind see [1]. Kneip first started the work of this direction by introducing ordered linear functionals (see [9]). For general results of alternating maximization (minimization) we refer to [2] and references therein.

The main point of the work is that the Alternating Method gives only a little gain, if any, in the complexity of optimum point computation for Linear Models (see [12]). Meanwhile, in

non-linear models the gain is pretty sensible. In non-linear models in most of the cases the closed form solution can not be obtained, in some cases even the numerical solutions of first order conditions could be very hard to implement in full parameter dimension. The technique known as alternating maximization (minimization) [3] helps in these situations and gives the estimation of parameter vector with adequate time complexity.

The model that we consider has the parameter  $\mathbf{v}$  which is of dimension  $p + q$ , where  $p$  is the dimension of *target* parameter and  $q$  is the dimension of *nuisance* parameter. Usually  $p$  is not large, because we also care about tractability and interpretability of our model, but  $q$  can be very large, although it is a *nuisance* parameter, we can not ignore it and exclude from considered model. Main problems with direct computations occur in high dimensions, i.e. in cases when  $p + q$  is large enough to make it impossible to invert matrix of sizes  $(p + q) \times (p + q)$ .

The alternating maximization procedure can be understood as a special case of the Expectation Maximization algorithm (EM algorithm). The EM algorithm is a popular algorithm first derived in [7]. Further on a number of modifications and extensions of this algorithm come into playground. In [7] it is also described how EM algorithm can be implement in different fields and give fruitful results. We refer to the [11] for the brief introduction to the development of EM algorithm. We restrict ourselves to citing the well-known convergence result (see [13]), which is still state of the art in most settings. Unfortunately, the result from [13] - as most convergence results on these iterative procedures - only ensures convergence to some set of local maximizers or fix-points of the procedure. In this work we consider one of the special cases where it is possible to show the actual convergence of the method.

The work has the following structure. In Section 2 contains the preliminaries about Generalized Linear Models (GLMs) and class of random variables called Exponential Family (EF). Section 3 contains the main results of the paper. Section 4 illustrates how the algorithm of alternating least squares (ALS) works and confirms above theoretical results using both real and simulated data.

## 2 Introduction to Generalized Linear Models (GLMs)

In this section we introduce to Generalized Linear Models (GLMs) from a slightly different angle. This section gives background for further research on Alternating Method for GLMs in Section 3.

Let  $Y_i$  be independent random vectors and regressors  $X_i \in \mathbb{R}^p$  and  $Y_i \sim P_i \in (\mathcal{P}_{\mathbf{v}})$ , which means that  $\exists v_i : P_i = P_{v_i}$ , where  $(\mathcal{P}_{\mathbf{v}})$  we assume to be an exponential family of distributions with canonical parameter. Exponential family will be discussed in details further in this section. Generalized linear model has the following representation  $Y_i \sim P_{v(X_i)}$ . In the case of Gaussian distribution we instead get the following form  $Y_i = v(X_i) + \varepsilon_i$  with arbitrary function  $v(\cdot)$ .

Function  $v(x)$  can be represented in the following way

$$v(x) = \sum_{j=1}^{\infty} \theta_j \psi_j(x).$$

Then, our linear parametric assumption is

$$v(x) = \sum_{j=1}^{p+q} \boldsymbol{\theta}_j \psi_j(x) = \sum_{j=1}^p \boldsymbol{\theta}_j \psi_j(x) + \sum_{j=p+1}^{p+q} \boldsymbol{\theta}_j \psi_j(x) \text{ for given basis } \psi_j(\cdot). \quad (2.1)$$

Denote  $\boldsymbol{\eta}_i = \boldsymbol{\theta}_{p+i}$  and column-vector  $\boldsymbol{\eta} = (\boldsymbol{\eta}_1, \dots, \boldsymbol{\eta}_q)^T \in \mathbb{R}^q$ . We will mostly discuss the case of finite  $p$  and  $q$ .

$$Y_i \sim P_{v_i}, \quad v_i = \Psi_i^T \boldsymbol{\theta} + \Phi_i^T \boldsymbol{\eta}, \quad (2.2)$$

where  $\Psi_i = (\psi_1(x), \dots, \psi_p(x))^T \in \mathbb{R}^p$ ,  $\Phi_i = (\psi_{p+1}(x), \dots, \psi_{p+q}(x))^T \in \mathbb{R}^q$  and  $\boldsymbol{\theta} \in \mathbb{R}^p$ ,  $\boldsymbol{\eta} \in \mathbb{R}^q$ . Log-likelihood

$$\log \frac{dP_{\boldsymbol{v}}}{d\mu_0^n}(\mathcal{Y}) = \sum_{i=1}^n (v_i Y_i - g(v_i)) = \sum_{i=1}^n (\Psi_i^T \boldsymbol{\theta} Y_i + \Phi_i^T \boldsymbol{\eta} Y_i - g(\Psi_i^T \boldsymbol{\theta} + \Phi_i^T \boldsymbol{\eta})),$$

which follows from the properties of exponential family, see (2.10). Function  $g(\cdot)$  is derived from (2.8) when the distribution is fixed. Equivalently, one writes

$$L(\boldsymbol{\theta}, \boldsymbol{\eta}) \stackrel{\text{def}}{=} S^T \boldsymbol{\theta} + R^T \boldsymbol{\eta} - A(\boldsymbol{\theta}, \boldsymbol{\eta}), \quad (2.3)$$

where

$$S \stackrel{\text{def}}{=} \sum_{i=1}^n Y_i \Psi_i \in \mathbb{R}^p, \quad R \stackrel{\text{def}}{=} \sum_{i=1}^n Y_i \Phi_i \in \mathbb{R}^q, \quad A(\boldsymbol{\theta}, \boldsymbol{\eta}) \stackrel{\text{def}}{=} \sum_{i=1}^n g(\Psi_i^T \boldsymbol{\theta} + \Phi_i^T \boldsymbol{\eta}).$$

Denote log-likelihood function as  $L(\boldsymbol{\theta}, \boldsymbol{\eta})$  or  $L(\boldsymbol{v})$ , where  $\boldsymbol{v} \stackrel{\text{def}}{=} (\boldsymbol{\theta}, \boldsymbol{\eta})$ . For the ease of representation sometimes we use  $\boldsymbol{v}$  instead of  $(\boldsymbol{\theta}, \boldsymbol{\eta})$  and vice versa.

For instance, in terms of  $\boldsymbol{v}$  log-likelihood can be written in the following way:

$$L(\boldsymbol{v}) = \mathcal{Y}^T \boldsymbol{v} - A(\boldsymbol{v}),$$

where  $\mathcal{Y} = \begin{pmatrix} S \\ R \end{pmatrix}^T \in \mathbb{R}^{p+q}$ , the concatenation of above defined  $S$  and  $R$ . Fisher information matrix is defined by  $\mathcal{D}^2 \stackrel{\text{def}}{=} -\nabla^2 \mathbb{E}L(\boldsymbol{v}^*) = F(\boldsymbol{v}^*)$ , where  $\boldsymbol{v} = \begin{pmatrix} \boldsymbol{\theta} \\ \boldsymbol{\eta} \end{pmatrix}^T \in \mathbb{R}^{p+q}$  and the Hessian matrix

$$\mathbb{F}(\boldsymbol{v}) = -\nabla^2 \mathbb{E}L(\boldsymbol{v}) = \begin{pmatrix} \mathbb{F}_{\boldsymbol{\theta}\boldsymbol{\theta}}(\boldsymbol{v}) & \mathbb{F}_{\boldsymbol{\theta}\boldsymbol{\eta}}(\boldsymbol{v}) \\ \mathbb{F}_{\boldsymbol{\eta}\boldsymbol{\theta}}(\boldsymbol{v}) & \mathbb{F}_{\boldsymbol{\eta}\boldsymbol{\eta}}(\boldsymbol{v}) \end{pmatrix}.$$

Further we will prove that function  $g(\cdot)$  is convex which gives us positively defined Fisher information matrix  $\mathbb{F}$ .

Also define score vector  $\nabla \stackrel{\text{def}}{=} \nabla L(\boldsymbol{v}^*)$  as well as standardized score vector  $\check{\xi}$  as follows:

$$\check{\xi} = \mathcal{D}^{-1} \nabla$$

Rewrite the definition (1.2) and (1.3) as follows

$$\tilde{\boldsymbol{v}} = (\tilde{\boldsymbol{\theta}}, \tilde{\boldsymbol{\eta}}) = \arg \max_{\boldsymbol{\theta}, \boldsymbol{\eta}} L(\boldsymbol{\theta}, \boldsymbol{\eta}), \quad \boldsymbol{v}^* = (\boldsymbol{\theta}^*, \boldsymbol{\eta}^*) = \arg \max_{\boldsymbol{\theta}, \boldsymbol{\eta}} \mathbb{E}L(\boldsymbol{\theta}, \boldsymbol{\eta}). \quad (2.4)$$

Parameters  $\tilde{\mathbf{v}} = (\tilde{\boldsymbol{\theta}}, \tilde{\boldsymbol{\eta}})$  are data-dependent and change by the change of  $Y$ , while  $\mathbf{v}^* = (\boldsymbol{\theta}^*, \boldsymbol{\eta}^*)$  is not random (i.e. data-dependent), but are unavailable under real world conditions. In real world we know nothing about distribution of  $Y$ , nevertheless we make parametric assumptions on  $Y$  and always keep in mind that parametric assumption is probably wrong. From the definition of  $\mathbf{v}^*$  follows that  $\nabla \mathbb{E}L(\mathbf{v}^*) = 0$  which yields

$$\mathbb{E} \begin{pmatrix} S \\ R \end{pmatrix}^T = \nabla A(\boldsymbol{\theta}^*, \boldsymbol{\eta}^*)$$

or equivalently

$$\mathbb{E}\mathcal{Y} = \nabla A(\mathbf{v}^*),$$

An important feature of exponential family is that the stochastic component  $\zeta(\boldsymbol{\theta}, \boldsymbol{\eta})$  of (log) likelihood function is linear in  $\boldsymbol{\theta}$  and  $\boldsymbol{\eta}$ . Put  $\varepsilon_i = Y_i - \mathbb{E}Y_i$  and  $\zeta = L - \mathbb{E}L$  then

$$\zeta(\boldsymbol{\theta}, \boldsymbol{\eta}) = (S^T - \mathbb{E}S^T)\boldsymbol{\theta} + (R^T - \mathbb{E}R^T)\boldsymbol{\eta} = \sum_{i=1}^n \varepsilon_i (\Psi_i^T \boldsymbol{\theta} + \Phi_i^T \boldsymbol{\eta}),$$

$$\nabla \zeta(\boldsymbol{\theta}, \boldsymbol{\eta}) = \begin{pmatrix} S - \mathbb{E}S & R - \mathbb{E}R \end{pmatrix} = \begin{pmatrix} \sum_{i=1}^n \varepsilon_i \Psi_i & \sum_{i=1}^n \varepsilon_i \Phi_i \end{pmatrix}.$$

Now consider the following elliptic set

$$\Omega_o(r) \stackrel{\text{def}}{=} \{\mathbf{v} : \|\mathcal{D}(\mathbf{v} - \mathbf{v}^*)\| \leq r\}. \quad (2.5)$$

Set  $\Omega_o(r)$  is called local vicinity of  $\mathbf{v}^*$  for  $\mathcal{D}^2 = -\nabla^2 \mathbb{E}L(\mathbf{v}^*) = \mathbb{F}(\mathbf{v}^*)$  and  $\mathcal{V}^2 \stackrel{\text{def}}{=} \text{Cov}(\nabla L(\mathbf{v}^*))$ .

Consider the covariance matrix in the block form:

$$\mathcal{V}^2 = \begin{pmatrix} V^2 & E \\ E^T & Q^2 \end{pmatrix}. \quad (2.6)$$

Now we want to understand the what these results yield for the estimator  $\tilde{\boldsymbol{\theta}}$  of the target parameter  $\boldsymbol{\theta}^*$ . Recall the block representation of the Fisher information matrix  $\mathbb{F}(\mathbf{v}^*) = -\nabla^2 \mathbb{E}L(\mathbf{v}^*)$ :

$$\mathbb{F}(\mathbf{v}) = \begin{pmatrix} \mathbb{F}_{\boldsymbol{\theta}\boldsymbol{\theta}}(\mathbf{v}) & \mathbb{F}_{\boldsymbol{\theta}\boldsymbol{\eta}}(\mathbf{v}) \\ \mathbb{F}_{\boldsymbol{\eta}\boldsymbol{\theta}}(\mathbf{v}) & \mathbb{F}_{\boldsymbol{\eta}\boldsymbol{\eta}}(\mathbf{v}) \end{pmatrix}. \quad (2.7)$$

For the central point  $\mathbf{v}^*$  the decomposition can be written in the following form

$$\mathcal{D}^2 = \mathbb{F}(\mathbf{v}^*) = \begin{pmatrix} D^2 & A \\ A^T & H^2 \end{pmatrix},$$

where  $D^2 = \mathbb{F}_{\boldsymbol{\theta}\boldsymbol{\theta}}(\mathbf{v}^*)$ ,  $A = \mathbb{F}_{\boldsymbol{\theta}\boldsymbol{\eta}}(\mathbf{v}^*)$  and  $H^2 = \mathbb{F}_{\boldsymbol{\eta}\boldsymbol{\eta}}(\mathbf{v}^*)$ . Decompose also the score vector  $\nabla \stackrel{\text{def}}{=} \nabla L(\mathbf{v}^*)$  as follows

$$\nabla = \begin{pmatrix} \nabla_{\boldsymbol{\theta}} \\ \nabla_{\boldsymbol{\eta}} \end{pmatrix}.$$

## 2.1 Exponential Family with Canonical Parameter (EFc)

In this part of our work we extend the scope of our modeling toolbox to accommodate a variety of additional data types, including counts and rates. We introduce the exponential family of distributions. It is worth mentioning that exponential family is pretty wide class of random variables and covers the vast majority of possible reasonable error distributions. It includes the Gaussian, binomial, Poisson, gamma, multinomial and many others. The simplest examples of distributions which are not from exponential family are Uniform distribution and Student's  $t$  distribution. The beauty and elegance of exponential family is the fact that the (log) likelihood function has a simple form and can be written explicitly.

In general we say that the random variable with probability density function  $f(\cdot)$  is from EFc, if the probability density function could be expressed in the following way

$$f(x|\boldsymbol{\eta}) = h(x) \exp\{\boldsymbol{\eta}^T T(x) - A(\boldsymbol{\eta})\} \quad (2.8)$$

where  $\boldsymbol{\eta}$  from l.h.s. is a given measure and the vector parameter  $\boldsymbol{\eta}$  from r.h.s. is referred to as the canonical parameter.  $T(X)$  is a *joint sufficient statistics* and we call function  $A(\cdot)$  as *link function*.

As stated above there is a huge number of famous distributions the probability density functions of which could be expressed in the form (2.8).

### 2.1.1 GLM with error distribution from EFc

Let  $Y$  be response vector of two sets of factors:  $\Psi$  and  $\Phi$ . Consider the following model

$$Y \sim P \in (\mathcal{P}_{\mathbf{v}})_{\mathbf{v} \in \mathbb{R}} \ll \mu_0, \quad (2.9)$$

where  $(\mathcal{P}_{\mathbf{v}})_{\mathbf{v} \in \mathbb{R}}$  is an exponential family with canonical parameter (EFc) and  $\mu_0$  is some dominating measure. Hence,

$$\log \frac{dP_{\mathbf{v}}}{d\mu_0} = y \cdot \mathbf{v} - g(\mathbf{v}) \text{ for some function } g. \quad (2.10)$$

**Lemma 2.1** (see []). *Let  $(\mathcal{P}_{\nu})$  be EFc. Then, it holds.*

$$\mathbb{E}_{\nu} Y = g'(\nu), \quad \text{Var}_{\nu}(Y) = \mathbb{E}_{\nu}[Y - g'(\nu)]^2 = g''(\nu) \quad (2.11)$$

Hence, function  $g(\cdot)$  is convex.

**Remark 2.1.** The above proof is provided for univariate functions  $g(\cdot)$ , but generally it is also true for multivariate functions. The only difference will be in the equation (??), instead of variance there will be a covariance matrix, which will be positive-defined.

## 2.2 Concentration

Recall that the following properties hold true for GLMs: the stochastic component  $\zeta(\boldsymbol{\theta}, \boldsymbol{\eta})$  is linear in both  $\boldsymbol{\theta}$  and  $\boldsymbol{\eta}$ , and the deterministic part  $\mathbb{E}L(\mathbf{v})$  is concave in  $\mathbf{v}$ . Consider the following elliptic set defined in (2.5). In [12] it is proved that there is a elliptic set  $\Omega_{\circ}(\mathbf{r})$  around oracle  $\mathbf{v}$  such that

$\Omega_o(\mathbf{r})$  contains the estimator  $\tilde{\mathbf{v}}$  with high probability. Further on we assume a sufficiently large value of  $\mathbf{x}$  to be fixed. It determines the level of overwhelming probability: a generic random set  $\Omega_0(\mathbf{x})$  is of dominating probability if

$$\mathbb{P}(\Omega_0(\mathbf{x})) \geq 1 - \mathcal{C}e^{-\mathbf{x}},$$

for some absolute constant  $\mathcal{C} > 0$ .

The value of  $\mathbf{x}$  may depend on sample size  $n$  and grows to infinity with  $n$ . The possible choices of  $\mathbf{x}$  are  $\mathbf{x} \asymp n^{1/2}$  and  $\mathbf{x} \asymp \log n$ , which would entail that  $\mathbb{P}(\Omega_0(\mathbf{x})) \geq 1 - \mathcal{C}/n$ . We only require that sequence  $\{x_n\}$  is not too large, more precisely,  $\mathbf{x} \leq \mathbf{x}_c \stackrel{\text{def}}{\asymp} n^{1/2}$ .

The way of construction and the exact probability of desired elliptic set is a very technical issue, hence omitted. All the results obtained below are considered to be true on the random set  $\Omega_o(\mathbf{x})$  of high probability.

### 2.3 Local Quadratic Approximation of the Log-likelihood

Recall that function  $L(\boldsymbol{\theta}, \boldsymbol{\eta})$  could be rewritten in terms of the whole parameter vector  $\mathbf{v}$  as  $L(\mathbf{v})$ . In this part of our work we will approximate our log-likelihood function  $L(\mathbf{v})$  with quadratic form in a local elliptic set of optimum  $\mathbf{v}^*$ . Put  $L(\mathbf{v}_1, \mathbf{v}_2) = L(\mathbf{v}_1) - L(\mathbf{v}_2)$  and recall that  $\tilde{\mathbf{v}}$  is data-dependent maximum of the likelihood function  $L(\cdot)$ , more precisely,

$$\tilde{\mathbf{v}} = \arg \max_{\mathbf{v}} L(\mathbf{v}), \quad \mathbf{v}^* = \arg \max_{\mathbf{v}} \mathbb{E}L(\mathbf{v}).$$

Then, using Taylor expansion for function  $L(\cdot)$  up to second order around the point  $\mathbf{v}^*$  yields

$$L(\mathbf{v}, \mathbf{v}^*) = \nabla L(\mathbf{v}^*)(\mathbf{v} - \mathbf{v}^*) - \frac{1}{2} \|\mathcal{D}^2(\mathbf{v} - \mathbf{v}^*)\|^2 + \alpha'(\mathbf{v}, \mathbf{v}^*), \quad (2.12)$$

where  $\alpha'(\cdot)$  is defined in (2.12).

Analogically, we approximate our function  $L(\mathbf{v})$  around  $\tilde{\mathbf{v}}$  together with the fact that  $\nabla L(\tilde{\mathbf{v}}) = 0$ , which gives us the following expression:

$$L(\mathbf{v}, \tilde{\mathbf{v}}) = -\frac{1}{2} \|\mathcal{D}^2(\mathbf{v} - \tilde{\mathbf{v}})\|^2 + \alpha(\mathbf{v}, \tilde{\mathbf{v}}). \quad (2.13)$$

**Remark 2.2.** In order to apply Taylor expansion we use the concentration property, which is crucial, because it helps us to claim that this expansion adequately represents the original function  $L(\cdot)$ . The idea comes from a well-known fact that every concave function is quadratic around its maximizer, hence it can be expanded in the Taylor series up to the second order without any significant losses.

## 3 Alternating Least Squares for GLMs

This chapter specifies the previously obtained general results in linear models to the case of a Generalized Linear Models (GLMs). The question is non-trivial because we can not obtain direct estimates and closed forms and we lose the quadraticity of log-likelihood function. Instead we will approximate our likelihood by a quadratic form in a set where the measure concentration is

observed. Generalized Linear Models are frequently used in many areas and applications including categorical data analysis, classification problem, Poisson and Binary regressions, statistical learning and density estimation. GLMs are popular in economics, because of its flexibility to various types of models. In economics sometimes linear regression is not enough to fully explain the phenomenon and deeply analyze the situation. Using GLMs partially solves this problem, but it makes problem from the semiparametric family, which is the payment for more reliable and accurate results. The restrictions of GLMs are not very strict, so that it can be easily adopted in various models.

Here we will discuss alternating method for the (quasi) likelihood function obtained in previous section. In general the alternating maximization (minimization) procedure is used in cases when the direct full dimension computations are not feasible or simply very difficult to implement.

Recall the log-likelihood function  $\mathcal{L}(\mathbf{v})$ , where vector  $\mathbf{v} = (\boldsymbol{\theta}, \boldsymbol{\eta})$  can be decomposed to the *target* parameter  $\boldsymbol{\theta}$  and *nuisance* parameter  $\boldsymbol{\eta}$ . Alternating maximization is an iterative procedure starting from initial value  $\mathbf{v}^\circ \in \mathbb{R}^{p+q}$  and updating iteratively in the way shown below

$$\begin{aligned}\tilde{\mathbf{v}}_{k,k} &\stackrel{\text{def}}{=} (\hat{\boldsymbol{\theta}}_k, \hat{\boldsymbol{\eta}}_k) = \left( \hat{\boldsymbol{\theta}}_k, \underset{\boldsymbol{\eta} \in \mathbb{R}^q}{\operatorname{argmax}} L(\hat{\boldsymbol{\theta}}_k, \boldsymbol{\eta}) \right), \\ \tilde{\mathbf{v}}_{k+1,k} &\stackrel{\text{def}}{=} (\hat{\boldsymbol{\theta}}_{k+1}, \hat{\boldsymbol{\eta}}_k) = \left( \underset{\boldsymbol{\theta} \in \mathbb{R}^p}{\operatorname{argmax}} L(\boldsymbol{\theta}, \hat{\boldsymbol{\eta}}_k), \hat{\boldsymbol{\eta}}_k \right).\end{aligned}\tag{3.1}$$

In this section we will try to answer to some natural questions that arise with the iterative procedure described above: Does the sequence  $(\hat{\boldsymbol{\theta}}_k)$  converge? And if the answer is yes, what is the convergence rate? What are the conditions under which the sequence actually converges to the global maximizer  $\tilde{\mathbf{v}}$ ?

### 3.1 Convergence to Likelihood Estimator

The following theorem is one of the main results about the convergence of alternating least squares for generalized linear models.

**Theorem 3.1.** *Suppose a model given by (2.2) and put  $\mathbf{v} = (\boldsymbol{\theta}, \boldsymbol{\eta}) \in \mathbb{R}^{p+q}$ .  $L(\mathbf{v})$  is defined in (2.3). Let  $\mathcal{D}^2 = -\nabla^2 \mathbb{E}L(\mathbf{v}^*)$  be the Hessian matrix for the log-likelihood function  $L(\mathbf{v})$  with the following block-matrix representation  $\begin{pmatrix} D^2 & A \\ A^T & H^2 \end{pmatrix}$  at point  $\tilde{\mathbf{v}}$ . Assume that  $\rho \stackrel{\text{def}}{=} \|D^{-1}AH^{-1}\|^2 < 1$  (here and afterwards the norm  $\|\cdot\|$  with matrix argument denotes as spectral norm of matrix) and the condition  $\|\mathcal{D}^{-1}\nabla^2 \mathbb{E}L(\mathbf{v})\mathcal{D}^{-1} - I_{p+q}\| \leq \delta(r) \leq \delta$  holds, where  $I_{p+q}$  is an identity matrix of size  $p+q$ .*

*Then the sequence of estimators  $(\hat{\boldsymbol{\theta}}_k)$  obtained by alternating least squares converges to  $\tilde{\boldsymbol{\theta}} = \Pi_{\boldsymbol{\theta}} \tilde{\mathbf{v}} \stackrel{\text{def}}{=} \Pi_{\boldsymbol{\theta}} \operatorname{argmax}_{\mathbf{v}} L(\mathbf{v})$ .  $\Pi_{\boldsymbol{\theta}}$  is the projection of full vector to its sub-vector  $\boldsymbol{\theta}$ .*

**Remark 3.1.** The notation  $\tilde{\mathbf{v}}_{k(+1),k}$  is used in cases when the results are true for both  $\tilde{\mathbf{v}}_{k,k}$  and  $\tilde{\mathbf{v}}_{k+1,k}$ .

*Proof.* Writing (2.13) in terms of  $\boldsymbol{\theta}$  and  $\boldsymbol{\eta}$  we get

$$L(\mathbf{v}, \tilde{\mathbf{v}}) =$$



$$-\frac{1}{2}\|D^2(\boldsymbol{\theta} - \tilde{\boldsymbol{\theta}})\|^2 - \frac{1}{2}\|H^2(\boldsymbol{\eta} - \tilde{\boldsymbol{\eta}})\|^2 - (\boldsymbol{\theta} - \tilde{\boldsymbol{\theta}})^T A(\boldsymbol{\eta} - \tilde{\boldsymbol{\eta}}) + \alpha(\mathbf{v}, \tilde{\mathbf{v}}). \quad (3.2)$$

Starting with an initial guess  $\boldsymbol{\theta}^\circ$  and using the method described in (3.1) we get

$$\hat{\boldsymbol{\eta}}_0 = \tilde{\boldsymbol{\eta}}(\boldsymbol{\theta}^\circ) =$$

$$\underset{\boldsymbol{\eta}}{\operatorname{argmin}} \left[ \frac{1}{2}\|D^2(\boldsymbol{\theta}^\circ - \tilde{\boldsymbol{\theta}})\|^2 + \frac{1}{2}\|H^2(\boldsymbol{\eta} - \tilde{\boldsymbol{\eta}})\|^2 + (\boldsymbol{\theta}^\circ - \tilde{\boldsymbol{\theta}})^T A(\boldsymbol{\eta} - \tilde{\boldsymbol{\eta}}) + \alpha(\mathbf{v}, \tilde{\mathbf{v}}) \right].$$

Hence, the first order condition gives us the following relationship:

$$H^2(\hat{\boldsymbol{\eta}}_0 - \tilde{\boldsymbol{\eta}}) = A^T(\tilde{\boldsymbol{\theta}} - \boldsymbol{\theta}^\circ) + \nabla_{\boldsymbol{\eta}}\alpha(\tilde{\mathbf{v}}_{0,0}, \tilde{\mathbf{v}}).$$

Analogically, the solution for  $\hat{\boldsymbol{\theta}}_1 \stackrel{\text{def}}{=} \tilde{\boldsymbol{\theta}}(\hat{\boldsymbol{\eta}}_0)$  has the following form

$$D^2(\hat{\boldsymbol{\theta}}_1 - \tilde{\boldsymbol{\theta}}) = A(\tilde{\boldsymbol{\eta}} - \hat{\boldsymbol{\eta}}_0) + \nabla_{\boldsymbol{\theta}}\alpha(\tilde{\mathbf{v}}_{1,0}, \tilde{\mathbf{v}}).$$

Then the iterative process of alternating maximization gives us the following recursive system of equations dependent on the initial guess:

$$\begin{cases} H^2(\hat{\boldsymbol{\eta}}_k - \tilde{\boldsymbol{\eta}}) = A^T(\tilde{\boldsymbol{\theta}} - \hat{\boldsymbol{\theta}}_k) + \nabla_{\boldsymbol{\eta}}\alpha(\tilde{\mathbf{v}}_{k,k}, \tilde{\mathbf{v}}) \\ D^2(\hat{\boldsymbol{\theta}}_{k+1} - \tilde{\boldsymbol{\theta}}) = A(\tilde{\boldsymbol{\eta}} - \hat{\boldsymbol{\eta}}_k) + \nabla_{\boldsymbol{\theta}}\alpha(\tilde{\mathbf{v}}_{k+1,k}, \tilde{\mathbf{v}}). \end{cases}$$

or, equivalently

$$\begin{cases} H^2(\hat{\boldsymbol{\eta}}_k - \tilde{\boldsymbol{\eta}}) = A^T(\tilde{\boldsymbol{\theta}} - \hat{\boldsymbol{\theta}}_k) + \nabla_{\boldsymbol{\eta}}\alpha(\tilde{\mathbf{v}}_{k,k}, \tilde{\mathbf{v}}) \\ D(\hat{\boldsymbol{\theta}}_{k+1} - \tilde{\boldsymbol{\theta}}) = D^{-1}A(\tilde{\boldsymbol{\eta}} - \hat{\boldsymbol{\eta}}_k) + D^{-1}\nabla_{\boldsymbol{\theta}}\alpha(\tilde{\mathbf{v}}_{k+1,k}, \tilde{\mathbf{v}}). \end{cases} \quad (3.3)$$

Then we express  $(\tilde{\boldsymbol{\eta}} - \hat{\boldsymbol{\eta}}_k)$  from the second equation of (3.3) and plug into the first one, which yields

$$D(\hat{\boldsymbol{\theta}}_{k+1} - \tilde{\boldsymbol{\theta}}) = D^{-1}AH^{-2}A^T(D^{-1}D)(\hat{\boldsymbol{\theta}}_k - \tilde{\boldsymbol{\theta}}) + D^{-1}[\nabla_{\boldsymbol{\theta}}\alpha(\mathbf{v}, \tilde{\mathbf{v}}) - AH^{-2}\nabla_{\boldsymbol{\eta}}\alpha(\mathbf{v}, \tilde{\mathbf{v}})],$$

then defining  $M_o \stackrel{\text{def}}{=} D^{-1}AH^{-2}A^TD^{-1}$  as follows and

$$\Xi(\tilde{\mathbf{v}}_{k(+1),k}) \stackrel{\text{def}}{=} D^{-1}[\nabla_{\boldsymbol{\theta}}\alpha(\tilde{\mathbf{v}}_{k+1,k}, \tilde{\mathbf{v}}) - AH^{-2}\nabla_{\boldsymbol{\eta}}\alpha(\tilde{\mathbf{v}}_{k,k}, \tilde{\mathbf{v}})]$$

gives us the following recursive formula

$$D(\hat{\boldsymbol{\theta}}_{k+1} - \tilde{\boldsymbol{\theta}}) = M_o \cdot D(\hat{\boldsymbol{\theta}}_k - \tilde{\boldsymbol{\theta}}) + \Xi(\tilde{\mathbf{v}}_{k(+1),k}). \quad (3.4)$$

Hence, summing up for all  $k$  starting from initial guess, taking the norm and using triangle inequality leads us to the following result:

$$\begin{aligned} \|D(\hat{\boldsymbol{\theta}}_{k+1} - \tilde{\boldsymbol{\theta}})\| &\leq \|M_o\| \cdot \|D(\hat{\boldsymbol{\theta}}_k - \tilde{\boldsymbol{\theta}})\| + \|\Xi(\tilde{\mathbf{v}}_{k(+1),k})\| \leq \|M_o\|^k \cdot \|D(\boldsymbol{\theta}^\circ - \tilde{\boldsymbol{\theta}})\| + \\ &\sum_{\ell=0}^{k-1} \|M_o\|^\ell \cdot \|\Xi(\tilde{\mathbf{v}}_{\ell(+1),\ell})\| = \|M_o\|^k \cdot \|D(\boldsymbol{\theta}^\circ - \tilde{\boldsymbol{\theta}})\| + \|\Xi(\tilde{\mathbf{v}}_{k(+1),k})\| \cdot \frac{1 - \|M_o\|^k}{1 - \|M_o\|}. \end{aligned}$$

Now using the assumption  $\rho \stackrel{\text{def}}{=} \|D^{-1}AH^{-1}\| < 1$  it is easy to see that the first term vanishes to zero when  $k$  goes to infinity.

We aim to show that the euclidean norm of  $\Xi(\tilde{\mathbf{v}}_{k(+1),k})$  adequately decreases with the rise of the number of iterations  $k$ .

Considering the term  $D^{-1}\nabla_{\boldsymbol{\theta}}\alpha(\tilde{\mathbf{v}}_{k+1,k},\tilde{\mathbf{v}})$  we refer to the Theorem C.1 in Andresen and Spokoiny (2013) which gives a vanishing in  $k$  bound for it. Note that using that theorem we can bound two remainder terms in  $\Xi(\tilde{\mathbf{v}}_{k(+1),k})$ , which will give the bound for  $\Xi(\tilde{\mathbf{v}}_{k(+1),k})$  as a whole.

Hence, using the Theorem from Andresen and Spokoiny (2013) and triangle inequality gives the claim about  $\Xi(\tilde{\mathbf{v}}_{k(+1),k})$ . More formally,

$$\|\Xi(\tilde{\mathbf{v}}_{k(+1),k})\| \leq \|D^{-1}\nabla_{\boldsymbol{\theta}}\alpha(\tilde{\mathbf{v}}_{k+1,k},\tilde{\mathbf{v}})\| + \|D^{-1}AH^{-2}\nabla_{\boldsymbol{\eta}}\alpha(\tilde{\mathbf{v}}_{k,k},\tilde{\mathbf{v}})\| \rightarrow 0 \text{ as } k \rightarrow \infty.$$

Combining all provided properties we derive that

$$\|D(\hat{\boldsymbol{\theta}}_{k+1} - \tilde{\boldsymbol{\theta}})\| \leq s_k \rightarrow 0 \text{ as } k \rightarrow \infty. \quad (3.5)$$

With the same argument it is clear how to derive the convergence property for *nuisance* parameter  $\boldsymbol{\eta}$ , which gives the claim of the theorem.  $\square$

**Remark 3.2.** The sequence  $s_k$  from (3.5) can be interpreted as the radius of the elliptic set around  $\tilde{\boldsymbol{\theta}}$  in which the estimator  $\hat{\boldsymbol{\theta}}_{k+1}$  derived by alternating least squares lies with high probability.

**Remark 3.3.** The above result presented in Theorem (3.1) tells us that we only need conditions  $\rho \stackrel{\text{def}}{=} \|M_{\circ}\| < 1$  and  $\|\mathcal{D}^{-1}\nabla^2\mathbb{E}L(\mathbf{v})\mathcal{D}^{-1} - I_{p+q}\| \leq \delta$  to claim that the sequence  $(\hat{\boldsymbol{\theta}}_k)$  derived by alternating maximization method will linearly converge to ML estimator  $\tilde{\boldsymbol{\theta}}$ , which is very often too difficult to compute directly.

### 3.2 Alternating Estimator

Above this section we have showed that the sequence of estimators obtained by using alternating least squares technique converges to the corresponding maximum likelihood estimator. It is a very strong and practically important result. As it has been discussed above, there are two main issues that make the problem non-trivial. The first one is impossibility of deriving closed form solution in most of cases and the second one is the high dimension of nuisance parameter which makes direct implementation of well-known Newton-Raphson method impossible. Alternating maximization technique overcame these issues and delivered estimators "close" to the maximum likelihood estimator.

This part is a bit technical and the only practical value is the statistical property of  $D(\hat{\boldsymbol{\theta}}_k - \boldsymbol{\theta}^*)$ .

Further in this section we will show that alternating estimator is also close to the true (unknown) parameter  $(\boldsymbol{\theta}^*, \boldsymbol{\eta}^*)$ . Recall

$$\mathbf{v}^* = (\boldsymbol{\theta}^*, \boldsymbol{\eta}^*) \stackrel{\text{def}}{=} \arg \max_{\mathbf{v}} \mathbb{E}L(\mathbf{v}).$$

The theorem below is known as a Fisher expansion and we formulate it in the framework of GLMs. Recalling the definitions given in the previous section we are now ready to formulate and prove Fisher expansion.

**Theorem 3.2.** Consider the conditions of Theorem (3.1) and define  $\check{\xi} \stackrel{\text{def}}{=} D^{-1}\check{\nabla}$ , where  $\check{\nabla} = \nabla_{\boldsymbol{\theta}} - AH^{-2}\nabla_{\boldsymbol{\eta}}$ .

Then we have

$$\|D(\tilde{\boldsymbol{\theta}}_k - \boldsymbol{\theta}^*) - \check{\xi}\| \rightarrow 0, \text{ as } k \rightarrow \infty.$$

*Proof.* The proof is based on the Theorem 1 proof's ideas with the only difference that the Taylor expansion is taken around point  $\mathbf{v}^*$ .

First, we use the first order conditions to obtain the following relations

$$\begin{cases} D(\tilde{\boldsymbol{\theta}}_k - \boldsymbol{\theta}^*) = D^{-1}\nabla_{\boldsymbol{\theta}}L(\mathbf{v}^*) - D^{-1}A(\tilde{\boldsymbol{\eta}}_k - \boldsymbol{\eta}^*) + D^{-1}\nabla_{\boldsymbol{\theta}}\alpha(\tilde{\mathbf{v}}_{k,k}, \mathbf{v}^*) \\ H(\tilde{\boldsymbol{\eta}}_k - \boldsymbol{\eta}^*) = H^{-1}\nabla_{\boldsymbol{\eta}}L(\mathbf{v}^*) - H^{-1}A^T(\tilde{\boldsymbol{\theta}}_{k-1} - \boldsymbol{\theta}^*) + H^{-1}\nabla_{\boldsymbol{\eta}}\alpha(\tilde{\mathbf{v}}_{k-1,k}, \mathbf{v}^*) \end{cases} \quad (3.6)$$

Again referring to Andreas and Spokoiny (2013) result we can bound the norm of  $D^{-1}\nabla_{\boldsymbol{\theta}}\alpha(\tilde{\mathbf{v}}_{k,k}, \mathbf{v}^*)$  and  $H^{-1}\nabla_{\boldsymbol{\eta}}\alpha(\tilde{\mathbf{v}}_{k-1,k}, \mathbf{v}^*)$ . From the system (3.6) one can easily derive

$$\begin{aligned} D(\tilde{\boldsymbol{\theta}}_k - \boldsymbol{\theta}^*) &= D^{-1}\nabla_{\boldsymbol{\theta}}L(\mathbf{v}^*) - D^{-1}[AH^{-2}\nabla_{\boldsymbol{\eta}}L(\mathbf{v}^*) - AH^{-2}A^T(\tilde{\boldsymbol{\theta}}_{k-1} - \boldsymbol{\theta}^*) + \\ &\quad AH^{-2}\nabla_{\boldsymbol{\eta}}\alpha(\tilde{\mathbf{v}}_{k-1,k}, \mathbf{v}^*)] + D^{-1}\nabla_{\boldsymbol{\theta}}\alpha(\tilde{\mathbf{v}}_{k,k}, \mathbf{v}^*) \implies \\ D(\tilde{\boldsymbol{\theta}}_k - \boldsymbol{\theta}^*) &= M_{\circ}D(\tilde{\boldsymbol{\theta}}_{k-1} - \boldsymbol{\theta}^*) + D^{-1}[\nabla_{\boldsymbol{\theta}}L(\mathbf{v}^*) - AH^{-2}\nabla_{\boldsymbol{\eta}}L(\mathbf{v}^*)] + \\ &\quad D^{-1}\{\nabla_{\boldsymbol{\theta}}\alpha(\tilde{\mathbf{v}}_{k,k}, \mathbf{v}^*) - AH^{-2}\nabla_{\boldsymbol{\eta}}\alpha(\tilde{\mathbf{v}}_{k-1,k}, \mathbf{v}^*)\}. \end{aligned} \quad (3.7)$$

Recalling the definitions from Section 3 we notice the term  $\check{\xi}$  as the second term of l.h.s. in (3.7), i.e. it can be rewritten in this way

$$D(\tilde{\boldsymbol{\theta}}_k - \boldsymbol{\theta}^*) - \check{\xi} = M_{\circ}D(\tilde{\boldsymbol{\theta}}_{k-1} - \boldsymbol{\theta}^*) + D^{-1}\{\nabla_{\boldsymbol{\theta}}\alpha(\tilde{\mathbf{v}}_{k,k}, \mathbf{v}^*) - AH^{-2}\nabla_{\boldsymbol{\eta}}\alpha(\tilde{\mathbf{v}}_{k-1,k}, \mathbf{v}^*)\}$$

Summing up for all  $k$  starting from the initial guess after taking the norm of both sides and using the assumption that  $\rho = \|M_{\circ}\| < 1$  leads us to the following result

$$\|D(\tilde{\boldsymbol{\theta}}_k - \boldsymbol{\theta}^*) - \check{\xi}\| \leq \mathcal{C}(\rho) \cdot D^{-1}\{\nabla_{\boldsymbol{\theta}}\alpha(\tilde{\mathbf{v}}_{k,k}, \mathbf{v}^*) - AH^{-2}\nabla_{\boldsymbol{\eta}}\alpha(\tilde{\mathbf{v}}_{k-1,k}, \mathbf{v}^*)\}, \quad (3.8)$$

where  $\mathcal{C}(\rho)$  is some constant depending on  $\rho$  and now, once again, we refer to the same Theorem C.1 from Andersen and Spokoiny (2013) to bound the term on the l.h.s. with the sequence which vanishes to zero when  $k$  goes to infinity.

Remember that, in order to obtain such bound for  $D^{-1}\nabla_{\boldsymbol{\theta}}\alpha(\tilde{\mathbf{v}}_{k,k}, \mathbf{v}^*)$  we need the condition  $\|\mathcal{D}^{-1}\nabla^2\mathbb{E}L(\mathbf{v})\mathcal{D}^{-1} - I_{p+q}\| \leq \delta$  for some constant  $\delta > 0$ .

Finally,

$$\|D(\tilde{\boldsymbol{\theta}}_k - \boldsymbol{\theta}^*) - \check{\xi}\| \leq \check{s}_k \rightarrow 0, \text{ as } k \rightarrow \infty. \quad (3.9)$$

□

Note that the expectation of random variable  $\check{\xi}$  is equal to zero and moreover  $\mathbb{V}ar(\check{\xi}) = D^{-1}V^2D^{-1}$ .

**Remark 3.4.** It is worth mentioning that in case of finite  $k$  Fisher expansion is still working, but the corresponding norm is bounded with the zero-vanishing sequence  $\check{s}_k$ , instead of exact zero.

Remark (3.5) comments about the distribution of  $\check{\xi}$  both asymptotically and in finite sample size.

**Remark 3.5.** Random variable  $\check{\xi} \in \mathbb{R}^p$  under the correct model specification has normal distribution, hence  $\|\check{\xi}\|^2$  has the distribution of  $\chi_p^2$  with  $p$  degrees of freedom. In case of model misspecification we can say nothing about the distribution of  $\|\check{\xi}\|^2$  in non-asymptotic framework, but asymptotically it still has  $\chi_p^2$  distribution with  $p$  degrees of freedom.

### 3.3 Importance of Initial Guess

The initial guess can play a crucial role in the convergence of alternating method and if we "succeed" with it then the gain is twofold. Firstly, we can come up with weaker conditions than in Theorem (3.1) and secondly, the number of iterations will be significantly less than in case of "bad" initial point. Good initial values of  $\theta^\circ$  mostly appeal to the first condition of Theorem (3.1), i.e.  $\rho \stackrel{\text{def}}{=} \|D^{-1}AH^{-1}\|^2 < 1$ .

Denote the vector space of eigenvectors corresponding to the eigenvalues greater than 1 of matrix  $M_\circ$  as  $\mathcal{V}$ , i.e.  $\mathcal{V} \stackrel{\text{def}}{=} \text{span}(v_1, v_2, \dots, v_l)$ , where

$$M_\circ v_i = \lambda_i v_i, \forall i \in \{1, \dots, l\} \quad \text{and} \quad |\lambda_i| \geq 1.$$

**Lemma 3.3.** *Then if the initial value of  $\theta^\circ$  is chosen such that it holds*

$$\mathbf{u} \perp \mathcal{V}, \tag{3.10}$$

where  $\mathbf{u} \stackrel{\text{def}}{=} D(\theta^\circ - \tilde{\theta})$  then we will have the convergence mentioned above.

The proof of this lemma is straightforward based on simple linear algebra. Nevertheless, we will briefly explain the idea. Note that the assumption  $\rho \stackrel{\text{def}}{=} \|D^{-1}AH^{-1}\|^2 < 1$  means that all eigenvalues are inside the unit circle so that the process will converge. It is true independent of the initial guess. Nevertheless, we can weaken this result by the "good" choice of initial values. If the matrix  $M_\circ$  has eigenvalues which are outside the unit circle then the initial guess could help to vanish them to zero being just orthogonal to the space of the corresponding eigenvectors.

It is also worth mentioning that the condition of Theorem (3.1)

$$\|\mathcal{D}^{-1}\nabla^2 \mathbb{E}L(\mathbf{v})\mathcal{D}^{-1} - I_{p+q}\| \leq \delta(r) \leq \delta$$

has nothing to do with initial values and only needs to bound the remainder term  $\Xi(\tilde{\mathbf{v}}_{k(+),k})$ . Thus, this condition can not be weakened depending on the initial point  $\theta^\circ$ .

## 4 Numerical Example

In this section, we introduce and implement an algorithm in case of a toy example. Suppose, there are  $n$  coins on the table and someone randomly chooses a coin and start flipping  $k$  times. For the ease of illustration and implementation we further assume that there are only two coins. Suppose the first coin has a probability of "head" outcome of  $p_1$  and the second coin —  $p_2$ . We introduce one *nuisance* parameter  $\pi$  which is the probability of choosing the first coin.

Assuming there are two coins on the table and one randomly chooses one of them 5 times and after each choice makes 10 flips. So, in total, we have 50 outcomes, 2 *target* parameters  $p_1$  and  $p_2$  and one nuisance parameter  $\pi$ .

We describe the algorithm on this simple example and provide the results of the code which computes the estimates of parameters in each step.

We have the following information:

Coin 2:  $H, T, H, T, T, H, T, H, H, T$ .

Coin 1:  $H, H, H, H, H, H, T, H, H, H$ .

Coin 1:  $H, T, H, T, H, H, H, H, H, H$ .

Coin 2:  $H, T, H, T, T, T, T, H, H, T$ .

Coin 1:  $H, H, H, T, H, H, T, H, H, T$ .

Clearly,

$$\tilde{p}_1 = \frac{24}{24+6} = 0.8 \quad \tilde{p}_2 = \frac{9}{9+11} = 0.45.$$

However, these information is not available for the algorithm. The algorithm starts with some initial values of probabilities, for example,  $p_1^o = 0.6$  and  $p_2^o = 0.5$ . Recall that in this example  $\boldsymbol{\theta}$  from the previous section is the vector  $(p_1, p_2)$  and  $\boldsymbol{\eta}$  is for  $\pi$ . The estimated probability of  $k$  out of 10 tosses of coin  $c \in \{1, 2\}$  yielding heads is

$$p_c(k) = C_k^{10} p_c^k (1 - p_c)^{10-k}. \quad (4.1)$$

Note that the binomial coefficient is the same for both coins, so it cancels out and only the ratio of the remaining factors determines the result. Using the initial guesses and formula (4.1) we come up with the following table

First iteration			
$\pi$	$1 - \pi$	Coin 1	Coin 2
0.45	0.55	$\approx 2.2H, 2.2T$	$\approx 2.8H, 2.8T$
0.80	0.20	$\approx 7.2H, 0.8T$	$\approx 1.8H, 0.2T$
0.73	0.27	$\approx 5.9H, 1.5T$	$\approx 2.1H, 0.5T$
0.35	0.65	$\approx 1.4H, 2.1T$	$\approx 2.6H, 3.9T$
0.65	0.35	$\approx 4.5H, 1.9T$	$\approx 2.5H, 1.1T$
		$\approx 21.3H, 8.6T$	$\approx 11.7H, 8.4T$

So, hence, in the next stage we get the following estimators of  $p_1$  and  $p_2$ :

$$\hat{p}_1^{(1)} = \frac{21.3}{21.3+8.6} = 0.71 \quad \hat{p}_2^{(1)} = \frac{11.7}{8.4+11.7} = 0.58. \quad (4.2)$$

The number of iterations until convergence depends on the tolerance level, which is the acceptable bias from the ML estimators. In other words, it defines the stopping criteria of the algorithm. At given initial values  $p_1^o = 0.6$  and  $p_2^o = 0.5$  of parameters  $p_1$  and  $p_2$  we get the following (see Table

Table 1: Convergence with random and non-random initial values

Iteration	Non-random		Random	
	$\hat{p}_1^{(i)}$	$\hat{p}_2^{(i)}$	$\hat{p}_1^{(i)}$	$\hat{p}_2^{(i)}$
1	0.600000000	0.500000000	0.935954410	0.0659086863
2	0.713012235	0.581339308	0.759177699	0.434881703
3	0.745292036	0.569255750	0.78052855	0.485752725
4	0.768098834	0.549535914	0.79025212	0.505705724
5	0.7831645	0.534617454	0.794205962	0.513867055
6	0.791055245	0.52628116	0.795774011	0.517227986
7	0.794532537	0.522390437	0.79639082	0.518613125
8	0.79592866	0.520729878	0.796632802	0.519183793
9	0.796465637	0.520047189	0.796727694	0.519418794
10	0.796668307	0.519770389	0.796764932	0.519515523
11	0.796744149	0.519658662	0.796779561	0.51955532
12	0.796772404	0.519613607	0.796785317	0.519571692
13	0.796782900	0.519595434		

1) sequence of estimated parameters. The table also contains the iterations of the algorithm with randomly chosen initial values from uniform distribution at certain seed.

In both cases we see that it converges to the ML estimator with high accuracy.

**Remark 4.1.** Note that the coins can be easily extended to any distribution. For example, the first part of observed data can be from one Gaussian distribution and the second part is from another Gaussian distribution. Formally,  $X_1, \dots, X_\ell \sim \mathcal{N}(\mu_1, \sigma_1^2)$  and  $X_{\ell+1}, \dots, X_n \sim \mathcal{N}(\mu_2, \sigma_2^2)$ . The idea is general and the same as in the case of coins (Bernoulli distribution) when we know that our data comes from  $n$  different normal distributions  $\mathcal{N}(\mu_i, \sigma_i^2)$ ,  $\forall i \in \{1, \dots, n\}$  and we introduce new parameters  $p_i$ , which are the probabilities of the corresponding distributions.

## References

- [1] Andresen A., (2015). Finite sample analysis of profile m-estimation in the single index model. *Electronic Journal of Statistics*, 9(2):2528–2641.
- [2] Andresen, A., Spokoiny, V., (2015) Two convergence results for an alternation maximization. arXiv: 1501.01525.
- [3] Bezdek, J., Hathaway, R. (2003). Convergence of Alternating Optimization. *Neural, Parallel & Pacific Computations* 11, 351 — 368.
- [4] Cavalier, L., Golubev, G. K., Picard, D., and Tsybakov, A. B. (2002). Oracle inequalities for inverse problems. *The Annals of Statistics*, 30(3): 843 — 874.

- [5] Chen, X. (2007). Large sample sieve estimation of semi-nonparametric models. *Handbook of Econometrics*, 6:55495632.
- [6] Chernozhukov, V., Chetverikov, D., and Kato, K. (2014). Anti-concentration and honest, adaptive confidence bands. *The Annals of Statistics*, 34(4):1653 — 1677.
- [7] A.P. Dempster, N.M. Laird, and D.B. Rubin (1977). Maximum likelihood from incomplete data via the EM algorithm. *Journal of the Royal Statistical Society, Series B*, 39:1 — 38.
- [8] Ibragimov, I.A. and R.Z. Khas'minskij (1981). "*Statistical estimation. Asymptotic theory.*" "New York - Heidelberg - Berlin: Springer-Verlag".
- [9] Kneip, A., (1994). Ordered linear smoothers. *The Annals of Statistics*, 22(2):835—866.
- [10] M. R. Kosorok. *Introduction to Empirical Processes and Semiparametric Inference* (2005). *Springer in Statistics*.
- [11] G. J. McLachlan and T. Krishnan (1997). *The EM Algorithm and Extensions*. Wiley, New York.
- [12] Spokoiny, V., Dickhaus, T., (2015). Basics of Modern Mathematical Statistics . *Springer Texts in Statistics*.
- [13] Wu, C. F. J. (1983). On the convergence properties of the EM algorithm. *Annals of Statistics*, 11:95 — 103.