

Chapter 3

Analytical and iterative factor modeling

Modeling a process which generated a given multivariate time series dataset for the purpose of learning motivates this thesis. In Chapter 2, the essential time series analysis tools that are needed in the modeling were presented; whereas in this chapter the elements of building the model itself will be discussed. Models with parameters that could be tuned to fit the statistical characteristics of the dataset at hand will be chosen; this tuning is called **parametric estimation**. It may be seen as a limitation because it warrants assumptions on the type of the data generation process involved. However, essential precautions will be taken by modeling on a dataset that is representative enough of the process. Moreover, in order to avoid any overfitting, learning methods which caution when wide deviations from the assumptions on the model and data are detected will be used.

The treatment of this chapter from the rest of the thesis has one main difference; here, any temporal correlation of the data samples in the given dataset is ignored. Yet, later on in the thesis, the parametric modeling techniques presented herein with the time series techniques of Chapter 2 will be utilized to achieve the thesis objectives.

In Section 3.1, a well-founded modeling strategy based on the **principle of maximum likelihood** will be introduced [96]. The principle assumes that the data has been generated by a known class of probability distributions whose parameters are to be estimated such that the **likelihood** of observing the data is maximized.

In Section 3.2, the concept of a **linear model** whose parameters are linear combinations of the data samples will be introduced. The derivation of the optimal parameters will be summarized by the **Gauss-Markov theorem**. The ideas of **unbiased** and **efficient** parameters defined there are desirable properties for any parametric model.

In Section 3.3, the **factor model** is presented. While the linear model of Section 3.2 utilizes some **measured variables** of the dataset to explain themselves or other measured variables, the factor model is remarkably different. The latter assumes existence of a fewer number of unmeasured **latent variables** responsible for generating all measured variables of a given dataset. The transformation from latent variables to measured variables is assumed to be non-random but unknown; this transformation will account for the covariations in the data. However, in the measured data, there will be deviations unexplained by such a generative model. Those deviations will be assumed unique to each of the measured variables and the variables that absorb these unique

deviations will be called **unique factors**. This characteristic of the factor model is actually facilitated by imposing a diagonal structure on the covariance matrix of the unique factors; whereas the latent variables are transformed such that they absorb the common variation of the measured variables. The transformed latent variables are, hence, called **common factors**.

Note that the parameters of the factor model are (i) **transformation matrix** of the latent variables to the measured variables and (ii) **variances of the unique factors**. The principle of maximum likelihood cannot yield these two sets of parameters independently. By assuming knowledge or guessing one of them, an estimate for the other parameter could be found.

In Section 3.4.1, the **principal factor model** first estimates the covariance matrix of the unique factors in order to estimate the transformation matrix; the reverse procedure is followed in **principal component factor model** of Section 3.4.2.

In Section 3.5, the **Expectation - Maximization (EM)** algorithm for maximum likelihood estimation of the factor model will be first narrated in an original manner. It is an iterative scheme established by [33]. The expression for **complete log-likelihood** of the measured variables as well as the latent variables are written out. However, its analytical tediousness in direct maximization is realized. This is overcome by probing its lower bound. It turns out that the local maximum of the lower bound is attained whenever the complete log-likelihood converges to the log-likelihood of the measured variables. Hence, starting with a set of guessed parameters, iteratively maximizing the complete log-likelihood converges towards the standard log-likelihood. Writing the lower bound of the complete log-likelihood scheme in an *a posteriori* expectation format and maximizing it for the optimal parameters is the crux of the EM algorithm.

In Sections 3.6 and 3.7, the scheme for using the EM algorithm for iteratively estimating factor model parameters is presented; it is partly along the lines of [14]. In doing so, the expression for the log-likelihood in the complete log-likelihood form is first written out; the latter is conducive for use with the algorithm. In the E-step of the algorithm presented in Section 3.7.1, the *a posteriori* mean and covariance of the latent variables are derived. In the M-step of Section 3.7.2, *a posteriori* expectation format of the log-likelihood is maximized; the parameters of the factor model, viz., transformation matrix of the common factors and covariance matrix of the unique factors, are thereby estimated.

3.1 Maximum likelihood model

This section starts by presenting some notions and usages that will help in explaining the characteristics of the data to be modeled. The primary assumption is that the data is a collection of samples of some relevant variables measured in an experiment; this collection will be denoted by \mathcal{D} and will be called simply the **dataset**. Let \mathcal{D} be constituted by n data samples written in the sequence $\mathcal{D} = \{y_l\}$, $l = 1, \dots, n$ and y_l be called the l -th **sample**.

In this chapter, unlike in Chapter 2, any sequential dependence of the value or occurrence of a data sample on next or any other sample is ignored. So there is no need to sort the data samples based on the time of data acquisition or any other criteria. But an index to identify the samples individually may be used.

Let the samples in \mathcal{D} be realizations of the sequence y_l of random variables. Let p^{y_l} denote density functions of their respective probability distributions. Throughout this

chapter, it is assumed that the samples encompassing \mathcal{D} adhere to the characteristic defined below [108]:

Definition 3.1. A dataset $\{y_l\}, l = 1, \dots, n$ due to its respective random variables y_l is said to be **independently and identically distributed or iid** if its joint probability density function is $\prod_{l=1}^n p^y(y_l)$, where y is a random variable such that $p^y(y_l) = p^{y_l}(y_l) \forall l = 1, \dots, n$.

Based on Definition 3.1, \mathcal{D} is interpreted as n realizations of a random variable y whose probability density function is p^y . Therefore, any realization y of y has a corresponding probability density $p^y(y)$. Extending $p^y(y)$ based on Definition 3.1 to \mathcal{D} leads to the joint probability density of the dataset, which is simply denoted by $p^y(\mathcal{D})$, and is given by

$$(3.1) \quad p^y(\mathcal{D}) \triangleq \prod_{l=1}^n p^y(y_l).$$

In order to maintain simplicity for the model which generated \mathcal{D} , a set of parameters θ will be introduced for the model. In the context of Definition 3.1, θ refers to the set of parameters of the probability density function of y . However, θ is unknown and has to be estimated. In this setup, $p^{y|\theta}(\mathcal{D} | \theta)$ will be termed as the likelihood of the dataset whilst the parameters θ are available. Note that the distribution is over y and the likelihood is a function of non-random θ .

At a set of parameters θ , due to (3.1), the likelihood of the dataset \mathcal{D} consisting of iid samples $y_l, l = 1, \dots, n$ factorizes as

$$(3.2) \quad p^{y|\theta}(\mathcal{D} | \theta) = \prod_{l=1}^n p^{y|\theta}(y_l | \theta).$$

In order to find an appropriate model for a given dataset, the intention is to utilize the following statistical methodology [96].

Definition 3.2. According to the **principle of maximum likelihood**, an optimal set of parameters $\hat{\theta}$ for the model corresponding to a dataset \mathcal{D} is the set of parameters θ for which the likelihood of \mathcal{D} is maximized, i.e.,

$$(3.3) \quad \hat{\theta} = \underset{\theta}{\operatorname{argmax}} p^{y|\theta}(\mathcal{D} | \theta).$$

A modification to (3.3) is made now: Since probability density is a non-negative function, the logarithm of the likelihood is maximized to arrive at the same solution for the optimal parameters as per Definition 3.2. Such an analysis of the exponential family of probability density functions will lead to substantial simplification [96]. Hence, (3.3) may be rewritten equivalently as

$$(3.4) \quad \hat{\theta} = \underset{\theta}{\operatorname{argmax}} L(\theta),$$

where

$$(3.5) \quad L(\theta) = \log_e p^{y|\theta}(\mathcal{D} | \theta).$$

3.2 Linear model

In Section 3.1, an appropriate estimate $\hat{\theta}$ of θ is considered as parameter for a model given the dataset \mathcal{D} . It is assumed that θ is a non-random quantity. Now consider the estimator Θ of θ , i.e., Θ is a random variable. It is hoped that Θ gives a reasonably good estimate of the true set of parameters θ given \mathcal{D} . Denoting mean of Θ by μ^Θ and variance by γ^Θ , the following properties indicate the quality of Θ ; refer §4.4 of [70] and §10.3 of [94]:

Property 3.1. An estimator Θ of θ is **unbiased** if $\mu^\Theta = \theta$.

Property 3.2. An estimator Θ of θ is **efficient** if $\Theta = \operatorname{argmin}_{\tilde{\Theta} \in \mathcal{C}} \gamma^{\tilde{\Theta}}$, where \mathcal{C} is the class of all unbiased estimators of θ .

A modeling strategy in Section 3.1 with a dataset constituted by iid samples was considered. Another popular parametric modeling paradigm called **linear model** involves treating a set of r -variate iid samples y_1, \dots, y_n as dependent on a set of q -variate iid samples x_1, \dots, x_n , where $n > q$. The simplest of linear models regresses x_l towards $y_l \forall l = 1, \dots, n$ through the relation

$$(3.6) \quad y_l = Wx_l + z_l,$$

where $W \in \mathbb{R}^{r \times q}$ is a linear function of $x_l \in \mathbb{R}^q$ and $z_l \in \mathbb{R}^r$ is the error in the regression [119, 109]. The model parameters θ discussed above refer to W here. Suppose the modeling errors z_l are realizations of the vector random variable z , then the measurements y_l may also be treated as realizations of the vector random variable y . The linear model may then be effectively written as

$$(3.7) \quad y = Wx + z,$$

for any $x \in \mathbb{R}^q$. By restricting the quality of the regression error z , the following theorem defines a popular linear model; for details one may refer to §6.2.1 of [62], §7.1 of [38] or §8.1 of [51] among plenty of references in the literature: Suppose each realization $y_l \in \mathbb{R}^r$ of the vector random variable y is related to $x_l \in \mathbb{R}^q$, $l = 1, \dots, n$ through (3.6) or (3.7) where z_l are due to zero mean uncorrelated Gaussian vector random variable z . According to the Gauss-Markov theorem, an efficient estimator of W is given by

$$(3.8) \quad \widehat{W} = [y_1 \cdots y_n]X'(XX')^{-1},$$

where $X = (x_1 \cdots x_n) \in \mathbb{R}^{q \times n}$ has $\operatorname{rank}(X) = q$. Then, the error estimate for the l -th sample becomes $\hat{z}_l = (\hat{z}_{1l} \cdots \hat{z}_{rl})' = y_l - \widehat{W}x_l \forall l = 1, \dots, n$. The unbiased estimator of the covariance matrix $\Gamma^z = \operatorname{diag}(\gamma^{z_1}, \dots, \gamma^{z_r})$ of z is given by

$$(3.9) \quad \gamma^{z_k} = \frac{1}{n-q} \sum_{l=1}^n \hat{z}_{kl}^2 \quad \forall k = 1, \dots, r.$$

3.3 Factor model

The linear model reviewed in Section 3.2 could be seen as the q -variables of the iid samples $x_l, l = 1, \dots, n$ together explaining the r -variables of each and every iid sample y_l , where both these sets of measured variables are available as part of the dataset \mathcal{D} . Suppose y_l and $x_l, l = 1, \dots, n$ are treated to be due to vector random variables y and x , respectively, so that y is the result of the transformation Wx , where $W \in \mathbb{R}^{r \times q}$. Then, the challenge is to explain y when x is unavailable or inaccessible in \mathcal{D} . One way to proceed is by assuming latent existence of the q -dimensional vector random variable x in generating the r -dimensional vector random variable y . In that context, y is named the set of **measured variables** and x the **latent variables**.

It is wished to pursue here a parametric model by involving the probability density function; this will help extract the statistical characteristics of the dataset in a finite number of parameters. Hence, if the probability density function of x is assumed known in the transformation $y = Wx$, then W serves as the parameter that needs to be estimated from the measured y .

However, the model $y = Wx$ is very restrictive because it assumes that any randomness in y is due to x whose characteristics are assumed. The model is relaxed by introducing an r -dimensional random variable z uncorrelated with x and designated to absorb all deviations in y that cannot be retained by

$$(3.10) \quad v \triangleq Wx.$$

Thus, the measured variables y are split into the common factors $v = Wx$ and unique factors z ; the following is such a model [80]:

Definition 3.3. A *factor model* is defined as

$$(3.11) \quad y = Wx + z,$$

where y and z are r -dimensional vector random variables, x is a q -dimensional vector random variable, $W \in \mathbb{R}^{r \times q}$ is a non-random transformation matrix, and

$$(3.12) \quad \mu^x = 0, \mu^y = 0, \mu^z = 0,$$

$$(3.13) \quad \Gamma^x = I_q,$$

$$(3.14) \quad \Gamma^z \text{ is diagonal, and}$$

$$(3.15) \quad \Gamma^{x,z} = 0.$$

Given a dataset of realizations of y , the parameters of the factor model that need to be estimated are W and the covariance matrix Γ^z of z . The factor model, in contrast to the linear model of (3.6), does not observe any realizations of x .

The following essential result for the moments of a function of a vector random variable is summarily provided; refer Chapter 6 of [35]: For $v = Wx$, $\mu^v = W\mu^x$ and $\Gamma^v = W\Gamma^xW'$. Therefore, applying (3.13) gives

$$(3.16) \quad \Gamma^v = WW'.$$

In the factor model, the condition (3.15) of zero correlation between x and z is crucial. Naturally, it leads to $\Gamma^{v,z} = 0$. Therefore, taking the second-order moments on both sides of (3.11) gives

$$(3.17) \quad \Gamma^y = \Gamma^v + \Gamma^z = WW' + \Gamma^z.$$

Due to (3.14), the r components of z are uncorrelated, or, all cross-covariances between the r components of y are inherited by only the covariance matrix $\Gamma^v = WW'$ of $v = Wx$ and not by Γ^z . This can be seen as each component of z inheriting only a part of the variance unique to its corresponding component in y . Hence, the components of z are called the **unique factors**. Since no part of the covariance of y common to all its components are held by z but instead by the transformation $v = Wx$, v is called the **common factors**.

Note that in the factor model, the r elements of W and r diagonal elements of Γ^z are to be estimated. Since (3.17) is just one equation with two unknowns, i.e., W and Γ^z , it cannot be solved uniquely; more conditions and assumptions may be placed to restrict possible solutions.

3.4 Maximum likelihood factor model

Well-known is the following assumption towards a proper solution of the factor model parameters, e.g., refer §3.5 of [62]: The measured variables follow a Gaussian distribution with parameters $\theta = \{\mu^y, \Gamma^y\}$, i.e.,

$$(3.18) \quad p^{y|\theta}(y | \theta) = \mathcal{N}(y | \mu^y, \Gamma^y),$$

as defined in (2.26).

Given samples y_l , $l = 1, \dots, n$ of the measured variables y , the principle of maximum likelihood as per Definition 3.2 could be used to estimate an optimal set of parameters according to (3.4). The maximum likelihood parameters $\hat{\mu}^y$ and $\hat{\Gamma}^y$ of the mean and covariance matrix of the Gaussian distribution in (3.18) are the sample mean and sample covariance matrix, respectively, i.e.,

$$(3.19) \quad \hat{\mu}^y = \frac{1}{n} \sum_{l=1}^n y_l,$$

$$(3.20) \quad \hat{\Gamma}^y = \frac{1}{n} \sum_{l=1}^n (y_l - \hat{\mu}^y)(y_l - \hat{\mu}^y)'$$

Then (3.20) may be substituted in (3.17) to get

$$(3.21) \quad \hat{\Gamma}^y = WW' + \Gamma^z.$$

However, it gives no clue regarding the maximum likelihood W and Γ^z , which are the parameters of interest to the factor model. In what follows, two relevant methods which derive appropriate solutions on the basis of the general maximum likelihood solution are briefly presented.

3.4.1 Principal factor model

One of the approaches to finding possible solutions to the maximum likelihood factor model of the r -dimensional measured variables y using q -dimensional common factors x and r -dimensional unique factors z as per Definition 3.3 starts with a good guess $\widehat{\Gamma}^z$ of Γ^z . The approach is called the principal factor model. One may refer to §10.2 of [54] or §6.3 of [46] to know how this guess could be made reliable; the details which are unnecessary for the objective of the present discussion are skipped. Substituting $\widehat{\Gamma}^z$ in (3.21) gives

$$(3.22) \quad \widehat{\Gamma}^y = WW' + \widehat{\Gamma}^z.$$

The problem then is to estimate a W such that $WW' = \widehat{\Gamma}^y - \widehat{\Gamma}^z$ subject to some quality criterion. Suppose columns u_1, \dots, u_r of $U \in \mathbb{R}^{r \times r}$ are the eigenvectors of $\widehat{\Gamma}^y - \widehat{\Gamma}^z$ whose corresponding eigenvalues $d_1^2 \geq \dots \geq d_r^2 > 0$ constitute the diagonal elements of a diagonal matrix D^2 from top-left to bottom-right. If the eigenvalue-eigenvector decomposition of $WW' = UD^2U'$ and the subscript $1 : q$ is used to refer the first q column indices of a matrix, then the optimal transformation matrix of the principal factor model is

$$(3.23) \quad \widehat{W} = U_{1:q}D_{1:q}.$$

If necessary, the estimation between $\widehat{\Gamma}^z$ and \widehat{W} may alternate iteratively.

3.4.2 Principal component factor model

Another approach, for which [103] is referred to, involves first estimating W and then the covariance matrix Γ^z of the unique factors. In order to estimate Γ^y and W , first, it has to be reminded that WW' is of rank q . Second, note that the relation (3.21) may be thought as WW' approximating the variance-covariance of the measured variables y as contained in Γ^y . There could be infinitely many ways WW' could approximate Γ^y and an approximation with respect to the Frobenius norm $\|\widehat{\Gamma}^y - \Gamma^y\|_F$ seems reasonable and standard practice; refer §10.2 of [54] and §2.12 of [103]. In that context, the following theorem is used; refer Lecture 5 of [120]:

Theorem 3.1. For full rank matrix $A \in \mathbb{C}^{r \times r}$ with eigenvectors u_1, \dots, u_r whose corresponding eigenvalues are $\alpha_1 \geq \dots \geq \alpha_r$, matrix $\widetilde{A} \in \mathbb{C}^{r \times r}$ with $\text{rank}(\widetilde{A}) = q < r$ defined as

$$(3.24) \quad \widetilde{A} = [u_1 \dots u_q] \text{diag}(\alpha_1, \dots, \alpha_q) [u_1 \dots u_q]^*$$

is such that

$$(3.25) \quad \|A - \widetilde{A}\|_F = \inf_{\substack{B \in \mathbb{C}^{r \times r} \\ \text{rank}(B)=q}} \|A - B\|_F = \alpha_{q+1}.$$

Due to Theorem 3.1, the optimal approximation of Γ^y using WW' in the Frobenius norm sense is declared as $\widehat{W}\widehat{W}' = E_{1:q}\Lambda_{1:q}^2E_{1:q}'$, where columns of E are eigenvectors of $\widehat{\Gamma}^y$ whose corresponding eigenvalues in decreasing order form the diagonal of Λ^2 .

Therefore, if the subscript $1 : q$ to refer to the first q column indices of a matrix, the estimate of W sought is given by

$$(3.26) \quad \widehat{W} = E_{1:q} \Lambda_{1:q}.$$

Since Γ^z ought to be diagonal due to uncorrelated z , taking into account (3.21), an approximate solution for Γ^z is

$$(3.27) \quad \widehat{\Gamma}^z \approx \widetilde{\text{diag}}(\widehat{\Gamma}^y - \widehat{W}\widehat{W}'),$$

where $\widetilde{\text{diag}}$ refers to setting the off-diagonal elements to zero.

3.5 EM algorithm

Now a presentation of the Expectation-Maximization (EM) algorithm is attempted; as stated in the introduction of this chapter, it is a popular iterative method for maximum likelihood estimation.

Note 3.1. *In this section, x is assumed a discrete and univariate random variable; this is to avoid any unnecessary analytical complications otherwise leading to equivalent conclusions. E.g. the summation over x has to be replaced by an integration for continuous x . And, a summation or integration across all dimensions of x is to be applied had x been a vector random variable.*

Note the following lemma (refer §4.5 of [49]):

Lemma 3.1. *If a random variable x is **marginalized** from its joint distribution $p^{x,y}$ with the random variable y , the result is the distribution of y , i.e.,*

$$(3.28) \quad p^y(y) = \sum_x p^{x,y}(x, y).$$

The definition of log-likelihood in (3.5) may be rewritten through Lemma 3.1 as

$$(3.29) \quad L(\theta) = \log_e \sum_x p^{y,x|\theta}(\mathcal{D}, x | \theta);$$

the maximization of this expression for the log-likelihood is intractable due to the summation inside the logarithm. In order to evade this situation, a dummy function $\eta(x)$ such that

$$(3.30) \quad \sum_x \eta(x) = 1; \quad \eta(x) > 0$$

is introduced and the complete log-likelihood

$$(3.31) \quad L(\theta, \eta) = \log_e \sum_x \eta(x) \frac{p^{y,x|\theta}(\mathcal{D}, x | \theta)}{\eta(x)}$$

is formed. The purpose in introducing $\eta(x)$ is to seek possibilities to maximize $L(\theta, \eta)$ in lieu of $L(\theta)$. In that pursuit, as shown in Section B.1.1, the logarithm may be brought inside the summation, i.e.,

$$(3.32) \quad L(\theta, \eta) \geq \sum_x \eta(x) \log_e \frac{p^{y, x | \theta}(\mathcal{D}, x | \theta)}{\eta(x)}.$$

Referring to Section B.1.2, it is possible to decompose the complete log-likelihood as

$$(3.33) \quad L(\theta, \eta) \geq L(\theta) + K(\theta, \eta),$$

where

$$(3.34) \quad K(\theta, \eta) = \sum_x \eta(x) \log_e \frac{p^{x | y, \Theta}(x | \mathcal{D}, \theta)}{\eta(x)}.$$

Now think of two iterative steps:

Step 1 – Find optimal η for a fixed θ : For a particular $\theta = \theta_i$, let $\hat{\eta}_i = \operatorname{argmax}_{\eta} L(\theta_i, \eta)$. Since local increase of $L(\theta, \eta)$ is guaranteed by locally maximizing its global lower bound $L(\theta) + K(\theta, \eta)$, $\hat{\eta}_i = \operatorname{argmax}_{\eta} K(\theta_i, \eta)$; one may refer [87] for more details. By differentiating (3.34) with respect to $\eta(x)$, it may be found that

$$(3.35) \quad \hat{\eta}_i = p^{x | y, \Theta}(x | \mathcal{D}, \theta_i).$$

However, $K(\theta_i, \hat{\eta}_i) = 0$ whereby

$$(3.36) \quad L(\theta_i, \hat{\eta}_i) = L(\theta_i).$$

Note that for a Gaussian density for y , the conditional probability for $\hat{\eta}_i$ in (3.35) is tractable.

Step 2 – Find optimal θ for a fixed η : Having found the locally optimal η for a fixed θ , the locally optimal θ for a fixed $\eta = \hat{\eta}_i$ is pursued. Based on (3.33) and (3.36), it may be written that

$$(3.37) \quad \theta_{i+1} = \operatorname{argmax}_{\theta} L(\theta, \hat{\eta}_i)$$

Note that (3.36) ensures that likelihood $L(\theta_i)$ is approached in every i -th iteration whenever $L(\theta_i, \hat{\eta}_i)$ is maximized to obtain the $i + 1$ -th estimate θ_{i+1} ; in other words, the iterations converge to a local maximum of $L(\theta_i)$.

E and M steps

The two steps arrived at above are now compiled. Suppose there is an initial guess θ_0 of θ . Then a local maximization of likelihood may be performed such that, in the i -th iteration, where $i = 1, 2, \dots$, has the two steps:

Step 1: estimate $\hat{\eta}_i$ based on θ_i , and

Step 2: locally maximize the likelihood to obtain θ_{i+1} .

These steps of an iteration become explicit if we note as shown in Section B.1.3 that

$$(3.38) \quad \theta_{i+1} = \operatorname{argmax}_{\theta_i} \mathbb{E}^{x|y,\theta} [\log_e p^{y,x|\theta}(\mathcal{D}, x | \theta_i)].$$

Hence, the i -th iteration involves:

Step 1 (Expectation): Evaluating the expectation $\mathbb{E}^{x|y,\theta} [\log_e p^{y,x|\theta}(\mathcal{D}, x | \theta_i)]$, and

Step 2 (Maximization): Maximizing $\mathbb{E}^{x|y,\theta} [\log_e p^{y,x|\theta}(\mathcal{D}, x | \theta_i)]$ locally with respect to θ_i .

3.6 Basic setup of EM algorithm for factor modeling

The difference between the linear model and the factor model is obvious and wide. While access to the samples x_n and y_n in (3.6) of the linear model is available, x in (3.11) is assumed inaccessible in the factor model. Hence, for the factor model, the conditional distribution $p^{y|x}$ of y given x is of interest.

Firstly, using the properties of the conditional distribution, e.g. refer §8 in Chapter 1 of [110], it could easily be shown for the factor model that $\mathbb{E}^{y|x}[y|x] = Wx + \mathbb{E}^{z|x}[z] = Wx$ because z is independent of x and has zero mean.

Secondly, the conditional variance $\Gamma^{y|x}$ is $\mathbb{E}^{y|x}[yy'|x] - (\mathbb{E}^{y|x}[y|x])(\mathbb{E}^{y|x}[y|x])'$. On expansion based on (3.11), $\Gamma^{y|x} = \mathbb{E}^{y|x} [(Wx + z)(Wx + z)'|x] - (Wx)(Wx)'$. Upon term-by-term expansion, due to the independence of z and x and since $\mathbb{E}^z[z] = 0$, the only surviving term will be $\mathbb{E}^{z|x}[zz'|x]$, which becomes $\mathbb{E}^z[zz'] = \Gamma^z$.

It is also well-known that the distribution of a Gaussian random vector conditioned on another is itself Gaussian; one may refer of §4.8 of [48] or Theorem 3.10.1 of [15] among many methods to verify it. Therefore, the Gaussian probability density of $y | x$ with parameters $\theta = \{W, \Gamma^z\}$ for the factor model may be written as

$$(3.39) \quad p^{y|x,\theta}(y | x, \theta) = \mathcal{N}(y | Wx, \Gamma^z).$$

Note that θ in $p^{y|x,\theta}$ refers to the availability of the set of parameters; the distribution is conditioned only on x . Based on the discussions in Section 3.1, the conditional probability density $p^{y|x,\theta}$ underpins the likelihood of the factor model. As with (3.1), the dataset \mathcal{D} is considered to consist of the iid samples y_l , $l = 1, \dots, n$ of y . The likelihood of the dataset is

$$(3.40) \quad p^{y|x,\theta}(\mathcal{D} | x, \theta) = \prod_{l=1}^n p^{y_l|x,\theta}(y_l | x, \theta).$$

Using Theorem B.1 known as Bayes theorem, $p^{y,x|\theta}(\mathcal{D}, x | \theta) = p^{y|x,\theta}(\mathcal{D} | x, \theta)p^x(x)$. If it is assumed that the distribution $p^x(x)$ to be independent of θ , then (3.38) of the EM-algorithm reduces to iteratively solving

$$(3.41) \quad \theta_{i+1} = \operatorname{argmax}_{\theta_i} \mathbb{E}^{x|y,\theta} [\log_e p^{y|x,\theta}(\mathcal{D} | x, \theta_i)].$$

From the above three equations, it may be written that

$$(3.42) \quad \theta_{i+1} = \underset{\theta_i}{\operatorname{argmax}} \mathbb{E}^{x|y,\theta} [f(\theta_i, x)],$$

$$(3.43) \quad \begin{aligned} f(\theta_i, x) &= \log_e p^{y|x,\theta}(\mathcal{D} | x, \theta_i) \\ &= \sum_{l=1}^n \log_e \mathcal{N}(y_l | W_i x, \Gamma_i^z), \end{aligned}$$

where the parameters

$$(3.44) \quad \theta_i \triangleq \{W_i, \Gamma_i^z\}$$

correspond to the i^{th} iteration.

3.7 Two steps of EM algorithm for factor modeling

In light of the discussion in Section 3.5 and the parameter update equations of (3.42), it may be stated that the i -th iteration of the factor model estimation alternates between:

1. **Expectation-step** Evaluate the expectation $\mathbb{E}^{x|y,\theta} [f(\theta_i, x)]$, and
2. **Maximization-step** Update $\theta_{i+1} \leftarrow \theta_i$ by maximizing $\mathbb{E}^{x|y,\theta} [f(\theta_i, x)]$ with respect to θ_i .

To proceed note that in (3.43) that $f(\theta_i, x) = -n \log_e(\det(\Gamma_i^z)) - \sum_{l=1}^n M(y_l, W_i x, \Gamma_i^z)$, where for any compatible vectors a, b and matrix C

$$(3.45) \quad M(a, b, C) = (a - b)' C^{-1} (a - b),$$

whose expansion gives $M(y_l, W_i x, \Gamma_i^z) = y_l' (\Gamma_i^z)^{-1} y_l - 2 y_l' (\Gamma_i^z)^{-1} W_i x + \operatorname{tr}((\Gamma_i^z)^{-1} W_i x x' W_i')$. Note the presence of terms with random variables x and $x x'$ in $M(y_l, W_i x, \Gamma_i^z)$. Therefore, the EM algorithm, as a result of this expansion of $M(y_l, W_i x, \Gamma_i^z)$, will involve alternating between:

1. **Expectation-step** Evaluate, for $l = 1, \dots, n$,

$$(3.46) \quad \begin{aligned} \langle x \rangle_{i,l} &\triangleq \mathbb{E}^{x|y,\theta} [x | y_l, \theta_i], \\ \langle x x' \rangle_{i,l} &\triangleq \mathbb{E}^{x|y,\theta} [x x' | y_l, \theta_i], \end{aligned}$$

where $\langle x \rangle_{i,l} \in \mathbb{R}^q$ and $\langle x x' \rangle_{i,l} \in \mathbb{R}^{q \times q}$, and

2. **Maximization-step** Update $\theta_{i+1} \leftarrow \theta_i$ by maximizing $f(\theta_i, x)$ with respect to θ_i , where x and $x x'$ are replaced by their corresponding *a posteriori* estimates, i.e., in (3.43)

$$(3.47) \quad \begin{aligned} x &\leftarrow \langle x \rangle_{i,l} \\ x x' &\leftarrow \langle x x' \rangle_{i,l}. \end{aligned}$$

The following analysis between (3.48) and (3.52) is inspired by [14].

3.7.1 E-step

Note that $\langle x \rangle_{i,l}$ is the mean of the Gaussian distribution $p^{x|y,\theta}(x | y_l, \theta_i)$, which is evaluated in Appendix B.2 to be

$$(3.48) \quad \begin{aligned} \langle x \rangle_{i,l} &= \Omega_i W_i' (\Gamma_i^z)^{-1} y_l, \\ \Omega_i &= (I_q + W_i' (\Gamma_i^z)^{-1} W_i)^{-1}, \end{aligned}$$

where $\Omega_i \in \mathbb{R}^{q \times q}$. From the classical relation of mean and covariance of any distribution, it is known that $\langle xx' \rangle_{i,l}$ is the sum of $\langle x \rangle_{i,l} \langle x' \rangle_{i,l}$ and covariance of $x | y_l, \theta_i$, i.e.,

$$(3.49) \quad \langle xx' \rangle_{i,l} = \langle x \rangle_{i,l} \langle x' \rangle_{i,l} + \Omega_i.$$

This completes the E-step of the EM Algorithm.

3.7.2 M-step

Towards the M-step of the EM algorithm, the substitutions in (3.47) give

$$(3.50) \quad \begin{aligned} E^{x|y,\theta}[f(\theta_i, x)] &= -n \log_e(\det(\Gamma_i^z)) - \sum_{l=1}^n \text{tr}((\Gamma_i^z)^{-1} W_i \langle xx' \rangle_{i,l} W_i') \\ &\quad - 2y_l' (\Gamma_i^z)^{-1} W_i \langle x \rangle_{i,l} + y_l' (\Gamma_i^z)^{-1} y_l. \end{aligned}$$

Now $E^{x|y,\theta}[f(\theta_i, x)]$ may be maximized to update the parameters W_i and Γ_i^z :

Update W_i : The problem that has to be solved is

$$(3.51) \quad \begin{aligned} W_{i+1} &= \underset{W_i}{\text{argmax}} E^{x|y,\theta}[f(\theta_i, x)] \\ &= \underset{W_i}{\text{arg}} \left(\frac{\partial}{\partial W_i} E^{x|y,\theta}[f(\theta_i, x)] = 0 \right). \end{aligned}$$

It is easy to see using matrix differentiation rules, e.g., refer [98], that

$$\frac{\partial}{\partial W_i} E^{x|y,\theta}[f(\theta_i, x)] = - \sum_{l=1}^n 2(\Gamma_i^z)^{-1} W_i \langle xx' \rangle_{i,l} - 2(\Gamma_i^z)^{-1} y_l \langle x' \rangle_{i,l},$$

which when equated to zero gives

$$(3.52) \quad W_{i+1} = \left(\sum_{l=1}^n y_l \langle x' \rangle_{i,l} \right) \left(\sum_{l=1}^n \langle xx' \rangle_{i,l} \right)^{-1}.$$

Update Γ_i^z : The access to the updated W_{i+1} is available and if

$$v_{i,l} = W_{i+1} \langle x \rangle_{i,l},$$

then $E^{x|y,\theta}[f(\theta_i, x)] = -n \log_e(\det(\Gamma_i^z)) - \sum_{l=1}^n M(y_l, v_{i,l}, \Gamma_i^z)$. Now, consider the update

$$(3.53) \quad \Gamma_{i+1}^z = \arg_{W_i} \left(\frac{\partial}{\partial W_i} E^{x|y,\theta}[f(\theta_i, x)] = 0 \right).$$

For $\Gamma_i^z = \text{diag}(\gamma_i^{z_1}, \dots, \gamma_i^{z_r})$ it can be seen that

$$E^{x|y,\theta}[f(\theta_i, x)] = -n \sum_{k=1}^r [\log_e(\gamma_i^{z_k}) + \frac{1}{\gamma_i^{z_k}} a_i^{z_k}]$$

where

$$a_i^{z_k} = \frac{1}{n} \sum_{l=1}^n (y_l - v_{i,l})^2.$$

Then, $\partial E^{x|y,\theta}[f(\theta_i, x)] / \partial \gamma_i^{z_k} = 0$ at

$$(3.54) \quad \gamma_{i+1}^{z_k} = a_i^{z_k}.$$

Factor model estimation via EM algorithm

Given the dataset, in Algorithm 2, the results of the analysis of the iterative parametric estimation of the factor model using the EM algorithm are summarized.

Algorithm 2: EM algorithm for the factor model

Input: $\mathcal{D} = \{y_l\}, l = 1, \dots, n$

Output: $\widehat{W}, \widehat{\Gamma}^z = \text{diag}(\widehat{\gamma}^{z_1}, \dots, \widehat{\gamma}^{z_r})$

initialize $i = 0$;

initialize randomly W_i, Γ_i^z ;

do

E-step:

for $l = 1$ *to* n **do**

 compute

$\langle x \rangle_{i,l}$ using (3.48);

$\langle xx' \rangle_{i,l}$ using (3.49);

end

M-step: update

W_{i+1} using (3.52);

$\gamma_{i+1}^{z_k} \forall k = 1, \dots, r$ using (3.54);

$i \leftarrow i + 1$;

$\epsilon \leftarrow E^{x|y,\theta}[f(\theta_i, x)] - E^{x|y,\theta}[f(\theta_{i-1}, x)]$ using (3.50);

while $\epsilon > 10^{-8}$ **and** $i < 20$;

$\widehat{W} \leftarrow W_i, \widehat{\gamma}^{z_k} \leftarrow \gamma_i^{z_k} \forall k = 1, \dots, r$;

A major drawback of the EM algorithm is the possibility that the estimation might get trapped in a local maximum of the log-likelihood and hence might require random restarts or other heuristic measures to be more certain regarding the estimates.

3.8 Summary

Two possibilities of modeling an r -dimensional measured vector random variable y were considered, viz., (i) the linear model where a measured variable $x \in \mathbb{R}^q, q < r$ is transformed to y and (ii) the factor model where a latent q -variate random variable x is transformed to y . Essentially, a factor model transforms a latent vector random variable of known probability distribution to a measured vector random variable of higher dimensionality that is perturbed by independent and uncorrelated noise. For the linear model, an efficient estimator of the transformation matrix was presented; whereas for the factor model there is no unique transformation. However, by restricting the variances unique to each of the measured variable, it is possible to estimate meaningful transformations. Thus, from a parametric modeling perspective, the transformation matrix and the unique variances are the parameters of the factor model.

In order to estimate the factor model parameters, two approaches based on the principle of maximum likelihood were discussed: The analytical estimation approach involves approximating the covariance structure of the measured variables using that of the transformed variables. For the iterative approach based on the EM algorithm, the log-likelihood function being lower bound by the *a posteriori* expectation of the logarithm of the joint probability density of the measured variables and the latent variables was exploited. Starting from guesses of the parameters, the EM algorithm maximizes the complete log-likelihood function of the latent variables and the measured variables by iteratively converging to the log-likelihood with every update of the parameters.