

Bayesian Techniques Applied in Counterfactual Analysis with High Dimensional Settings and Time Varying Properties ^{*}

Yaohan Chen

School of Economics, Singapore Management University

Last version: June 12, 2020

This version: December 15, 2020

Abstract

Bayesian methodologies are rapidly becoming tools in machine learning analysis. Among all these conventional and recently developed Bayesian analysis tools, Expectation Maximization (EM) algorithms plays the pivotal role in modern Bayesian analysis and so is its extension Variational Expectation and Maximization (VEM) algorithm. Both of these methodologies have been widely used in the application of modern Bayesian analysis in machine learning such as textual analysis established on Latent Dirichlet Allocation (LDA), which is usually referred to Topic Models. This note mainly discusses the modern technical tools and analytical framework for application of modern Bayesian analysis along with the associated properties with specific focus on Topic Models and how EM algorithm is set up from Variational Bayesian perspective.

Keywords: Bayesian; Machine Learning; Asset pricing; Dynamic Covariance Structure

^{*} I am grateful for the discussion kindly shared literatures with Associate Professor Tao Zeng from Zhejiang University. All errors are my own. You may contact me at yaohan.chen.2017@phdecons.smu.edu.sg.

1 Introduction

Variational Bayes methods and Empirical Bayes methods (Armstrong et al., 2020) have recently gained much attention in academia especially in machine learning literature since under many circumstances shrinkage has to be imposed to improve the both the in-sample fitting and out-of-sample prediction accuracy. Another major reason for the increasing popularity of this analysis framework within academia is that currently there is no general way to account for uncertainties without imposing many parametric assumptions, this is actually an issue that has been pointed in some recent literature like Hansen (2016, p.116). Bayesian analysis instead offers a flexible and philosophically coherent framework based on posterior analysis for learning with data rich in hierarchical structure. However in most practical application settings, the posterior is in general intractable and consequently some modern sampling algorithms like MCMC have to be adopted for approximating the target posterior based on the sampling. One key feature of sampling methods like MCMC is the intrinsically required heavy block-loop structure and hence very often computational efficiency is deteriorated in large-scale applications. By comparison, Variational Bayes (VB) method as an alternative for sampling algorithm like MCMC has been empirically proved to be able to circumvent the computational inefficiency often encountered in sampling algorithm while keeping desired approximation as much as possible (Braun and McAuliffe, 2010; Blei, Kucukelbir, and McAuliffe, 2017). VB actually is a specific methodology originated from the more extensive variational approximation analytical framework which can be at least traced back to (Jordan et al., 1999; Winn and Bishop, 2005). One of the major reasons for the documented computation efficiency gained from VB, which is one of the main attractive properties justifying its increasing popularity in large-scale application as well, is the corresponding feature in formulating the posterior approximation as a specific optimization problem. Briefly speaking, the rough idea underlying the core of VB is to find the closest distribution to the exact posterior over some family of distributions (Guo, Wang, Fan, Broderick, and Dunson, 2016). Recently there is also some work justifying the ground theory for VB (Wang and Blei, 2019; Alquier and Ridgway, 2020), where they have formerly defined the VB methodology and related theoretical properties of VB approximation (specifically all these work discuss the frequentist concentration of VB approximation for posterior under different divergence concepts). In fact VB is an extensive framework such that many classical algorithms like EM algorithm are possible to be included in the discussion within this framework. As for the development of modern empirical methods in Economics and Finance, Armstrong, Kolesár, and Plagborg-Møller (2020) extends empirical Bayes confidence interval of Morris (1983) in a robust way free of the prior distribution imposed on means.

Another reason why Variational Bayes method has attained much attention in economics literature is the development of economic models with intrinsic requirement to incorporate stochastic variation of parameters since Cooley and Prescott (1976). In recent decades, accompanied with the fast development of data science and machine learning techniques, econometric models with time-varying parameters and high-dimensional data-structure are broadly acknowledged in academia

as an important topic, there has been a large amount of recent literature established on this framework.

Another area where such a kind of analytical framework is possibly to be applied is counterfactual analysis with synthetic controls in high-dimensional settings. Seminal work can be traced back to Abadie and Gardeazabal (2003) and Abadie et al. (2010). Recently there is a nice and comprehensive review work on this framework by Abadie (2020). Ever since the synthetic control method was proposed in Abadie and Gardeazabal (2003) for the first time, it has received widespread attention both in economic and machine learning literature. The reason why counterfactual analysis is an important tool for social science is that under many circumstances causality is of the most importance to be discussed in social science researches (for instance there is growing literature in finance discussing the publication effect of anomaly-based trading strategies (McLean and Pontiff, 2016; Chen and Zimmermann, 2020b; Pelger and Xiong, 2020)) and usually the casual statements with respect to treatment (intervention) depends on the construction of counterfactuals, which are typically the unobserved outcomes that would have been if a unit had not been treated. Routinely, estimation of counterfactual effects of this kind is based on the aforementioned synthetic control method that makes the desired estimation upon the estimation of outcomes of similar group of individuals not affected by the intervention. This intuitive idea driving synthetic method is similar to the comparative studies widely adopted in other social science where the key idea is laid upon that the effect of intervention can be inferred by comparing the evolution of the target (outcome) variables of interest between the unit exposed to treatment and a group of units that are similar to the unit but were not affected by the treatment. A inherent key factor for ensuring the validity comparative analysis is the affinity between unit exposed to treatment and the units not exposed to treatment, which is however not that easy to be justified in practice. There is an emerging literature on applying machine learning methods in recovering the unobserved counterfactual effects from synthetic controls (Doudchenko and Imbens, 2017; Chernozhukov et al., 2020a,b; Athey et al., 2020).

One alternative for counterfactual construction discussed in Carvalho et al. (2018), the so-called ArCo method, is essentially a two-step procedure: In the first step, the data before the occurrence of the intervention is used to estimate a multivariate time-series regression model relating the variables in the treated (dependent variables) with the variables belonging to the untreated peers (explanatory variables); In the second step, the counterfactual is constructed by extrapolating the estimated model with data after intervention. Finally the estimated effect is the time-series average of the difference between the data actual data and counterfactual.

A natural follow-up question would be whether it is possible to account for the potential time-varying properties.

2 Basic Setup and Assumptions

2.1 Notation

All random variables are defined in a Probability Space $(\Omega, \mathcal{F}, \mathbb{P})$. Random variables are defined by upper case letter and the associated realization is defined by lower case letter such that $X(\omega) = x$. Matrices and vectors are written in bold letters \mathbf{X} . Hence each entry of the realized \mathbf{X} , i.e., x_{ij} , denotes the i -th realization (row) of the j -th random variable (column). Sets are denoted by calligraphic upper case \mathcal{X} and the associated cardinality is denoted by $|\mathcal{X}|$. For the symbol of norm, we use $\|\cdot\|$ to denote the generic (semi)norm; $\|\cdot\|_q$ and $\|\cdot\|_{\mathcal{L}^q}$ are adopted to denote ℓ^q and \mathcal{L}^q norms respectively for $q \in [1, \infty)$. Specifically these two norms are defined separately for the cases when $d > 1$ and $d = 1$ when the vector $\mathbf{X} = (X_1, \dots, X_d)^\top$ degenerates to a scalar (For this case we just use X to denote the degenerated scalar random variable). That is

$$\begin{aligned}\|\mathbf{X}\|_q &:= \left(\sum_{i=1}^d |X_i|^q \right)^{1/q} \\ \|X\|_{\mathcal{L}^q} &:= (\mathbb{E}|X|^q)^{1/q}\end{aligned}$$

We further introduce other notation for norm such that

$$\begin{aligned}\|\mathbf{X}\|_\infty &:= \max_{i \leq d} |X_i| \text{ if } \mathbf{X} \text{ is } d \times 1 \text{ vector} \\ \|\mathbf{X}\|_{\max} &:= \max_{i \leq m, j \leq n} |X_{i,j}| \text{ if } \mathbf{X} \text{ is } m \times n \text{ matrix}\end{aligned}$$

Furthermore, we introduce the following notation (ℓ^0 norm)

$$\|\mathbf{X}\|_0 := |\{i : X_i \neq 0\}|$$

and quadratic form associated with d -dimensional matrix \mathbf{M} takes the following form

$$\|\mathbf{X}\|_{\mathbf{M}}^2 := \mathbf{X}^\top \mathbf{M} \mathbf{X}$$

And we also adopt $\text{diag}(\mathbf{X})$ to denote the diagonal matrix whose diagonal elements are extracted from the diagonal of \mathbf{X} ; $\mathbf{1}(A)$ represents an indicator function on the event A , i.e., $\mathbf{1}(A) = 1$ if A is true or $\mathbf{1}(A) = 0$, otherwise. [Yes you can build a VAR model for the target variable and the covariates using the approach by Koop and Korobilis. Then the estimated model together with updated values for the covariates can be used to estimate counterfactual values of the target variable.](#) In the following discussion, I may use term “intervened” or “treated” interchangeably with the same meaning that intervention has taken place at a specific period.

Definition 1 *The probability density function for the random matrix \mathbf{X} ($n \times p$) that follows the*

matrix normal distribution $\mathcal{MN}(\mathbf{M}, \mathbf{U}, \mathbf{V})$ has the form:

$$p(\mathbf{X} | \mathbf{M}, \mathbf{U}, \mathbf{V}) = \frac{\exp\left(-\frac{1}{2} \text{tr}\left[\mathbf{V}^{-1}(\mathbf{X} - \mathbf{M})^\top \mathbf{U}^{-1}(\mathbf{X} - \mathbf{M})\right]\right)}{(2\pi)^{np/2} |\mathbf{V}|^{n/2} |\mathbf{U}|^{p/2}} \quad (1)$$

The matrix norm is related to the multivariate normal distribution in the following way

$$\mathbf{X} \sim \mathcal{MN}(\mathbf{M}, \mathbf{U}, \mathbf{V}) \quad (2)$$

if and only if

$$\text{vec}(\mathbf{X}) \sim \mathcal{N}(\text{vec}(\mathbf{M}), \mathbf{V} \otimes \mathbf{U})$$

where \mathbf{V} is $p \times p$ matrix and \mathbf{U} is $n \times n$ matrix respectively.

2.2 Counterfactual analysis and synthetic control

The procedure implemented in counterfactual analysis with synthetic controls is briefly discussed in this subsection. We assume the target variable is a scalar denoted by Y_t , which is exposed to a treatment (intervention) that occurred at $t = T_0 + 1$. Counterfactual effect is then estimated from the covariates \mathbf{X} (\mathbf{X} is generally a $P \times J$ matrix with the j -th column indicating the j -th unit, either treated or not; and for each column as a $P \times 1$ column vector refers to the vector collecting the covariates associated with each unit), i.e., characteristics of peers that are assumed to be unaffected by the intervention. Realization of real valued random vector is observed, which is denoted by

$$\mathbf{Z}_t = (Z_{1t}, \dots, Z_{J+1,t})^\top = (\mathbf{Z}_{0t}^\top : \mathbf{Z}_{1t}^\top)^\top$$

where we just introduce notation \mathbf{Z}_{0t}^\top and \mathbf{Z}_{1t}^\top as two sub column vectors of \mathbf{Z}_t to emphasize that the observed random vector generally can be separated into treated group (\mathbf{Z}_{1t}^\top accordingly) and untreated group (\mathbf{Z}_{0t}^\top accordingly). $\mathcal{D}_t \in \{0, 1\}$ is introduced as a sequence of binary variable indicating the periods when treatment (intervention) was in place, that is

$$\mathcal{D}_t = \begin{cases} 1 & \text{if } t > T_0 \\ 0 & \text{otherwise} \end{cases} \quad (3)$$

For each period t , the observed Z_{it} could be either from $Z_{it}^{(1)}$ which is the potential outcome when the unit i is exposed to the intervention or from $Z_{it}^{(0)}$ which is the potential outcome when the unit i is not exposed to the intervention. Hence we are able to jointly express $Z_{it} \in \mathbf{Z}_{1t}$ as

$$Z_{it} = \mathcal{D}_t Z_{it}^{(1)} + (1 - \mathcal{D}_t) Z_{it}^{(0)} \quad (4)$$

Without loss of generality, we may focus on the case when there is a single variable contained treated group and denote this scalar $Y_t = Z_{1t}$ as target variable; All the remained J units from donor

pool (peers assumed to be not affected by intervention) are collected in \mathbf{Z}_{0t} , that is

$$\mathbf{Z}_{0t} = (Z_{2t}, \dots, Z_{J+1,t})^\top$$

Given this restricted case when Z_{1t} is the target variable potentially exposed to intervention, we are interested in evaluating what would Z_{1t} have been like had there been no intervention, i.e. we are interested in obtaining the evaluation of $Z_{1t}^{(0)}$ by approximating it using the following functional form ¹

$$Z_{1t}^0 = \mathcal{M}\left(\mathbf{Z}_{0t}^{(0)}; \boldsymbol{\theta}_0\right) + V_t, \quad t = 1, \dots, T \quad (5)$$

where

$$\mathcal{M} : \mathcal{Z} \times \Theta \rightarrow \mathbf{R}, \quad \mathcal{Z} \subseteq \mathbf{R}^{n-1}$$

and then the estimated evaluation of $Z_{1t}^{(0)}$ is

$$\widehat{Z}_{1t}^{(0)} = \mathcal{M}\left(\mathbf{Z}_{0t}^{(0)}; \widehat{\boldsymbol{\theta}}_{T_0}\right)$$

The remained question of potential research interest is what exactly the functional form of $\mathcal{M}(\cdot)$ we are able to apply to approximate the target unobserved counterfactual effect. Machine learning method as one alternative of the statistician toolkit is widely known for its ability at this kind of objective, for instance the recently developed machine learning algorithms taking the advantage of Bayesian techniques featuring the tree structure (Chipman et al., 1998; Denison et al., 1998; Chipman et al., 2010; Ročková and van der Pas, 2019; Ročková, 2019; Ročková and Saha, 2019). Alternatively it is also worthwhile making an attempt to apply the framework proposed by Koop and Korobilis (2020) to model $\mathcal{M}(\cdot)$, since this would potentially bring the advantage of accounting for the time-varying properties of $\mathcal{M}(\cdot)$.

Alternatively, we may make our analysis restricted to linear framework and accordingly some theoretic properties are relatively easy to establish, hence this is by far the framework within which most of the discussions are made in literature. In particular, within linear framework, the general objective (Abadie, 2020) is to uncover the optimized \mathbf{W}^* such that for a given \mathbf{V} the following objective function is minimized

$$\|\mathbf{X}_1 - \mathbf{X}_0 \mathbf{W}\| = \left(\sum_{p=1}^P v_p (X_{p1} - w_2 X_{p2} - \dots - w_{J+1} X_{p,J+1})^2 \right)^{1/2} \quad (6)$$

¹ As argued in Abadie (2020), one of the key feature of synthetic control method is the implicit assumption that the combination of units in the donor pool may approximate the characteristics of the the affected unit substantially better than any unaffected unit alone, and this serves as one motivating reason why we want to introduce the general functional form of units in donor pool.

where

$$\begin{aligned}\mathbf{W} &= (w_2, \dots, w_J)^\top \\ \mathbf{V} &= (v_1, \dots, v_k)^\top\end{aligned}$$

and we abbreviate the timing subscript for both \mathbf{X}_{1t} and \mathbf{X}_{0t} with focus on the cross-sectional dimension. Then the functional form of counterfactual effect within this framework is

$$\widehat{Z}_{1t}^{(0)} = \mathcal{M}(\mathbf{Z}_{0t}; \widehat{\boldsymbol{\theta}}_{T_0}) = \sum_{j=2}^{J+1} w_j^* Z_{jt} \quad (7)$$

One of the key assumption justifying our discussion is

$$Z_{1t} \perp\!\!\!\perp \mathcal{D}_s \text{ for all } t, s.$$

Once we are able to back out $Z_{1t}^{(0)}$ as precise as possible, then conceptually the following difference is used as the proxy for evaluating counterfactual effects,

$$\mathcal{H}_0 : \delta_t := Z_{1t}^{(1)} - Z_{1t}^{(0)} \text{ for } t > T_0 \quad (8)$$

and corresponding estimation for δ_t along with it is given as

$$\widehat{\delta}_t = Z_{1t} - \widehat{Z}_{1t}^{(0)}$$

and finally the evaluated difference due to the intervention is

$$\widehat{\Delta}_T = \frac{1}{T - T_0} \sum_{t=T_0+1}^T \widehat{\delta}_t$$

and the associated $\frac{1}{T - T_0} \sum_{t=T_0+1}^T \delta_t$. In this note I also attach the replicated figures separately in [Figure 4\(a\)](#) and [Figure 4\(b\)](#) as following. As emphasized in Abadie (2020), synthetic control method is essentially a weighted average algorithm using the weighted average of characteristics of units from donor pool (i.e. sometimes it is referred to fitting) to approximate the unobserved counterfactual effect of target variable by restricting the weights on probability simplex. As we have always emphasized, the goal of counterfactual analysis or synthetic control is to approximate the trajectory that would have been observed for $Y_t = Z_{1t}$ at $t > T_0$ in the absence of intervention.

To justify the intrinsic advantage of synthetic control method in comparison to other methodologies, it would be useful to compare synthetic control method with linear regression framework. The way we adhere to using linear regression is augmenting the characteristic space of units from treated group and untreated group by adding a row of ones to \mathbf{X}_1 and \mathbf{X}_0 respectively, denoted by $\overline{\mathbf{X}}_1$ and $\overline{\mathbf{X}}_0$. Further more, let \mathbf{Y}_0 as the $(T - T_0) \times J$ matrix denote the observed effects for J units

from untreated group. The regression of \mathbf{Y}_0 on $\bar{\mathbf{X}}_0$ yields the estimated counterfactual effect based on regression as

$$\widehat{\mathbf{B}}^\top \bar{\mathbf{X}}_1 \quad (9)$$

where $\widehat{\mathbf{B}}$ collects the regression coefficients in the following form

$$\widehat{\mathbf{B}} = \left(\bar{\mathbf{X}}_0 \bar{\mathbf{X}}_0^\top \right)^{-1} \bar{\mathbf{X}}_0 \mathbf{Y}_0^\top$$

Hence this implies that the estimated counterfactual effect is

$$\underbrace{\mathbf{Y}_0 \bar{\mathbf{X}}_1^\top \left(\bar{\mathbf{X}}_0 \bar{\mathbf{X}}_0^\top \right)^{-1} \bar{\mathbf{X}}_1}_{\mathbf{w}^{\text{reg}}} \quad (10)$$

Remark 2.1 *For the restricted case where there is a single unit in the treated group, both \mathbf{X}_1 and $\bar{\mathbf{X}}_1$ are $P \times 1$ column vector, which implies the final estimated counterfactual effect collected in $\mathbf{Y}_0 \bar{\mathbf{X}}_1^\top \left(\bar{\mathbf{X}}_0 \bar{\mathbf{X}}_0^\top \right)^{-1} \bar{\mathbf{X}}_1$ is a $(T - T_0) \times 1$ column vector with each entry indicating the estimated counterfactual effect at each period after intervention.*

Remark 2.2 *One prominent feature of the prevalent synthetic analysis framework is the underlying implicit assumption that researchers are informed about [how the treated and untreated group is divided](#), but it seems not always the case in practice. Actually once the intervention is imposed, which units are exposed to the intervention generally is not known for researchers. Or in other words, for the target variable that is supposed to be contained in treated group, the corresponding treated effect on the target variable is time-varying (after the period when the intervention is imposed). Hence the counterfactual effect like a “phantom” should be treated as latent variable in general, and recently there is an inspiring work by [Bojinov and Shepard \(2019\)](#). In this work they have emphasized the timing scheme of counterfactual effect, which is demonstrated in the following adapted figure*

[Place [Figure 1](#) about here]

The basic idea is as following: it is obvious from the figure that for a realized sequence of target variable(s), the status whether it is (they are) affected by the intervention (treatment) is totally determined by the unobserved path

$$\mathbf{W}_{1:T} = (W_1, \dots, W_T)$$

where each element W_t ($1 \leq t \leq T$) as the binary variable indicating whether the associated target variable is exposed to intervention or not. W_t ($1 \leq t \leq T$) is also referred to treatment path in literature.

An straightforward extension in Bayesian framework is by considering the Dirichlet allocation or the variational bayes method widely adopted in Bayesian literature. Since all these techniques

automatically satisfy the inherent requirement of synthetic control method that the average weights have to be added up to 1. Moreover, as has been widely discussed in statistical literature, Dirichlet allocation intrinsically features the sparsity desired in the application of synthetic method. A specific discussion within Bayesian framework is discussed as following. Moreover, it is a natural idea to regard the unobserved counterfactual of units from treated group as latent variables within Bayesian framework.

Generally, there are two kinds of available information (data) for backing out the counterfactual effect: (i) the inherent trend information of target variable(s) contained in treated group, which is potentially able to be extracted by using (V)AR model; (ii) information contained in covariates collected in donor pool (untreated group), information of this kind in general features the dynamic of covariates. Finally the evaluated counterfactual effect can be written as the linear combination of the counterfactual effects estimated from these two kinds of information. Formally in mathematics,

$$\begin{cases} \mathbf{z}_{1t} = A\mathbf{z}_{1t-1} + \boldsymbol{\epsilon}_t & (11) \\ \mathbf{z}_{0t} = A\mathbf{x}_{0t} + \boldsymbol{\epsilon}_t & (12) \end{cases}$$

3 Latent Dirichlet Allocation

3.1 Basic concepts

As proposed in Blei et al. (2003), some basic concepts about text documents are introduced as following

- A **word** refers to a item drawn from vocabulary indexed by $\{1, \dots, V\}$, each word is represented by a V -vector such that if this word belongs to the v -th term then w is specifically a vector $w^v = 1$ and $w^u = 0$ if $u \neq v$.
- A **document** refers to a sequence of N words collected in $\mathbf{w} = (w_1, \dots, w_N)$, where w_n is the n -th word in sequence.
- A **corpus** is a collection of M documents denoted by $D = \{\mathbf{w}_1, \dots, \mathbf{w}_M\}$.

3.2 Estimation

Basically, Latent Dirichlet Allocation (LDA) assumes the following generating process,

Step 1 The term distribution $\boldsymbol{\beta}$ is determined for each topic by

$$\boldsymbol{\beta} \sim \text{Dir}(\boldsymbol{\delta})$$

Step 2 Choose $\theta \sim \text{Dir}(\boldsymbol{\alpha})$.

Step 3 For each of N words w_n :

3a) Choose a topic $z_n \sim \text{Multinomial}(\theta)$.

3b) Choose a word w_n from $p(w_n | z_n, \beta)$, a multinomial probability conditional on the topic z_n .

where β refers to a word probabilities matrix of dimension $k \times V$ with,

$$\beta_{ij} = P(w^j = 1 | z^i = 1).$$

Specifically more in details, a k -dimensional Dirichlet random variable can take values in the $(k-1)$ simplex (a k vector lies in the $(k-1)$ -simplex if $\theta_i \geq 0$ and $\sum_{i=1}^k \theta_i = 1$), and has the following probability density function over its simplex

$$\begin{aligned} p(\theta | \alpha) &= \frac{\Gamma(\sum_{i=1}^k \alpha_i)}{\prod_{i=1}^k \Gamma(\alpha_i)} \theta_1^{\alpha_1-1} \dots \theta_k^{\alpha_k-1} \\ &= \frac{\Gamma(\alpha k)}{\Gamma(\alpha)^k} \prod_{i=1}^k \theta_i^{\alpha-1} \quad (\text{A special case in which } \alpha_1 = \dots = \alpha_k) \end{aligned}$$

Remark 3.1 *How Dirichlet distribution is updated based on data observations. Suppose that random variables are drawn from multinomial distribution (thus sampling from $1, \dots, k$ with repetition) with the sampling probability drawn from the special case of Dirichlet distribution mentioned in the context. There are totally N -times samplings such that*

$$x_1 + \dots + x_k = N \text{ and } 1 \leq i \leq k$$

where x_i represents how many times k is selected. The joint density is

$$\begin{aligned} p(\mathbf{x}, \theta | \alpha) &= \frac{N!}{x_1! \dots x_k!} \prod_{i=1}^k \theta_i^{x_i} \cdot \frac{\Gamma(\alpha k)}{\Gamma(\alpha)^k} \prod_{i=1}^k \theta_i^{\alpha-1} \\ &= \frac{N!}{x_1! \dots x_k!} \frac{\Gamma(\alpha k)}{\Gamma(\alpha)^k} \cdot \prod_{i=1}^k \theta_i^{x_i + \alpha - 1} \end{aligned} \quad (13)$$

where (13) is of the functional form of the density function of Dirichlet distribution, hence it is easy to obtain the marginalization of the distribution of data \mathbf{x} by integrating out θ , which yields

$$p(\mathbf{x} | \alpha) = \frac{N!}{x_1! \dots x_k!}$$

Consequently according to Bayes rule, we have the posterior update about θ based on data observations as

$$p(\theta | \mathbf{x}, \alpha) = \frac{p(\mathbf{x}, \theta | \alpha)}{p(\mathbf{x} | \alpha)} = \frac{\Gamma(\alpha k)}{\Gamma(\alpha)^k} \cdot \prod_{i=1}^k \theta_i^{x_i + \alpha - 1} = \text{Dir}(x_1 + \alpha, \dots, x_k + \alpha) \quad (14)$$

Discussion following conventional style in machine learning and natural language processing

literature generally adopts conventional methodologies for processing unstructured information extracted from text, the corresponding analysis of which relies much on counting the occurrences of words. These methodologies however are justified by the implicit assumption that the order of words collected in documents (filings) is not related to the information to be recovered. By contrast, topic models featuring the application of LDA (Latent Dirichlet Allocation) provides a relatively successful extension of the conventional methodologies in natural language processing by assuming that the words collected in documents are implicitly in connection with several different topics (to some extent, topics or themes are regarded as the latent variables which are not observed and required to be estimated). To the end, topic models are essentially the modelling framework through which the aggregate information contained in documents (each document is regarded as a sequence of words in order) is extracted. Generally there are two alternatives for estimating this model, VEM (Variational Expectation Maximization) and Gibbs Sampling, which are to be briefly discussed in the following.

VEM (Variational Expectation Maximization) is introduced for estimation with the likelihood function specified as following

$$\begin{aligned} \ell(\alpha, \beta) &= \log(p(\mathbf{w} \mid \alpha, \beta)) \\ &= \log \int \left\{ \sum_{\mathbf{z}} \left[\prod_{i=1}^N p(\mathbf{w}_i \mid z_i, \beta) p(z_i \mid \theta) \right] \right\} p(\theta \mid \alpha) d\theta. \end{aligned}$$

where $\mathbf{z} = (z_i)_{i=1, \dots, N}$ includes all the combinations of assigning N words into k topics. VEM is specifically a combination of Variational inference and E-M algorithms. More details about the Variational inference refer to (Mclachlan and Krishnan, 2008; Dempster et al., 1977; Wainwright and Jordan, 2008). Since here entries of β are directly regarded as parameters to be estimated if VEM is applied, hyperparameter δ is no more the parameter of interest.

Gibbs sampling can be done as well for approximating $p(\mathbf{z} \mid \mathbf{w})$, which essentially corresponds to the posterior probability of topic conditional on the observed words. α and δ are fixed at suggested values with $\alpha = 50/k$ and $\delta = 0.1$. Theoretical work done by Griffiths and Steyvers (2004); Phan et al. (2008) has justified that the whole loop of Gibbs sampling is compositionally based on sampling from the following distribution Fama and French (1993, 1996, 2015)

$$p(z_i = K \mid \mathbf{w}, z_{-i}) \propto \frac{n_{-i,K}^{(j)} + \delta}{n_{-i,K}^{(\cdot)} + V\delta} \frac{n_{-i,K}^{(d_i)} + \alpha}{n_{-i,\cdot}^{(d_i)} + k\alpha} \quad (15)$$

- where $n_{-i,K}^{(j)}$ denotes how often the j -th term collected in vocabulary is assigned to topic K without the i -th word.
- d_i indexes the document to which word w_i belongs.
- Dot \cdot implies that the summation over this index is implemented.

The term distribution and topic distribution is updated as following

$$\hat{\beta}_K^{(j)} = \frac{n_K^{(j)} + \delta}{n_K^{(\cdot)} + V\delta} \quad \hat{\theta}_K^{(d)} = \frac{n_K^{(d)} + \alpha}{n^{(d)} + k\alpha} \quad (16)$$

for $j = 1, \dots, V$ and $d = 1, \dots, D$.

4 How to define Disaster Measure using LDA

The output from LDA generally gives us the key words extracted from text. The remained question is about how to match the extracted keywords (topics) with the existing dictionary indicating whether the associated contents are about disasters.

5 Notes for Variational Methods in Bayesian Analysis

5.1 EM Algorithm

Rather than sampling as in MCMC, EM algorithm initially proposed by (Dempster et al., 1977) chooses the target density function from a family of approximate density functions \mathcal{Q} such that,

$$q^*(\mathbf{z}) = \arg \min_{q(\mathbf{z}) \in \mathcal{Q}} \text{KL}(q(\mathbf{z}) \| p(\mathbf{z} | \mathbf{x})) \quad (17)$$

Here we just denote $q(\mathbf{z}) \in \mathcal{Q}$ without specifying any parameters for it, however sometimes for the implementability concern in practice we usually consider the parametrized family \mathcal{Q} such that each element contained in \mathcal{Q} takes the form of $q_\phi(\mathbf{z})$. Actually specification of the functional form of $q(\mathbf{z})$ depends on how much the flexibility to be taken into account for analysis and for the scenario when we make the assumption that $q(\mathbf{z})$ is factorizable, it can be demonstrated how optimization is set up in the discussion about Variational Expectation Maximization algorithm contained in the next subsection.

The most popular estimation technique adopted in statistical literature is MLE (Maximum Likelihood Estimation) where parameters are estimated based on the observed data by solving the following optimization

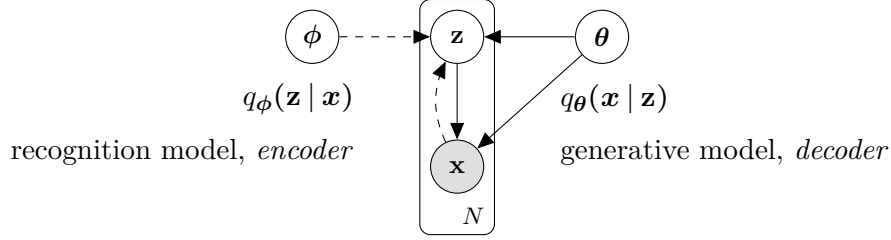
$$\hat{\theta}_{\text{ML}} = \arg \max_{\theta} p(\mathbf{x}; \theta). \quad (18)$$

Instead of being directly connected with observed data \mathbf{x} through data generating process, it is practical and sometimes for convenience concern to introduce latent variables \mathbf{z} such that the

likelihood function is actually the marginalization of joint likelihood as following ²

$$p(\mathbf{x}; \boldsymbol{\theta}) = \int p(\mathbf{x}, \mathbf{z}; \boldsymbol{\theta}) d\mathbf{z} = \int p(\mathbf{x} | \mathbf{z}; \boldsymbol{\theta}) p(\mathbf{z}; \boldsymbol{\theta}) d\mathbf{z} \quad (19)$$

It would also be helpful to summarize the general ideas for the framework in the following figure where latent variables are introduced, which is usually the case where Variational Bayes method plays the role.



Proposition 5.1 *Under some regular conditions,*

$$p(\mathbf{x}; \boldsymbol{\theta}) = F(q, \boldsymbol{\theta}) + \text{KL}(q \| p) \quad (20)$$

with

$$F(q, \boldsymbol{\theta}) = \int q(\mathbf{z}) \ln \left(\frac{p(\mathbf{x}, \mathbf{z} | \boldsymbol{\theta})}{q(\mathbf{z})} \right) d\mathbf{z} \quad (21)$$

Proof.

$$\begin{aligned} \ln p(\mathbf{x}; \boldsymbol{\theta}) &= \int q(\mathbf{z}) \ln \frac{p(\mathbf{x}; \boldsymbol{\theta}) p(\mathbf{z} | \mathbf{x}; \boldsymbol{\theta})}{p(\mathbf{z} | \mathbf{x}; \boldsymbol{\theta})} d\mathbf{z} \\ &= \int q(\mathbf{z}) \ln \frac{p(\mathbf{z}, \mathbf{x}; \boldsymbol{\theta})}{p(\mathbf{z} | \mathbf{x}; \boldsymbol{\theta})} d\mathbf{z} \\ &= \int q(\mathbf{z}) \ln \frac{p(\mathbf{z}, \mathbf{x}; \boldsymbol{\theta}) q(\mathbf{z})}{p(\mathbf{z} | \mathbf{x}; \boldsymbol{\theta}) q(\mathbf{z})} d\mathbf{z} \\ &= \int q(\mathbf{z}) \ln \left(\frac{p(\mathbf{x}, \mathbf{z} | \boldsymbol{\theta})}{q(\mathbf{z})} \right) d\mathbf{z} + \int q(\mathbf{z}) \ln \frac{q(\mathbf{z})}{p(\mathbf{z} | \mathbf{x}; \boldsymbol{\theta})} d\mathbf{z} \end{aligned}$$

It is known that Kullback-Leibler divergence is not negative and hence $F(q, \boldsymbol{\theta})$ refers to the Evidence Lower Bound (ELBO). □

² Actually an inspiring perspective on this would be regarding \mathbf{z} as local latent variables, one per observation, and $\boldsymbol{\theta}$ as (a d -dimensional vector and usually d is less than the sample size denoted by n) as global latent variables, which intrinsically does not change with sample size.

Remark 5.1 Alternatively, we may note that similar trick applies to the posterior with contingency on data. Usually this is for the case where Bayesian techniques apply when our target is to obtain the posterior mode, that is to maximize $\ln p(\boldsymbol{\theta}; \mathbf{x})$. Note that analogously we have

$$\begin{aligned}
\ln p(\boldsymbol{\theta}; \mathbf{x}) &= \int q(\mathbf{z}) \ln \frac{p(\boldsymbol{\theta}; \mathbf{x}) p(\mathbf{z} | \boldsymbol{\theta}; \mathbf{x})}{p(\mathbf{z} | \boldsymbol{\theta}; \mathbf{x})} d\mathbf{z} \\
&= \int q(\mathbf{z}) \ln \frac{p(\mathbf{z}, \boldsymbol{\theta}; \mathbf{x})}{p(\mathbf{z} | \boldsymbol{\theta}; \mathbf{x})} d\mathbf{z} \\
&= \int q(\mathbf{z}) \ln \frac{p(\mathbf{z}, \boldsymbol{\theta}; \mathbf{x}) q(\mathbf{z})}{p(\mathbf{z} | \boldsymbol{\theta}; \mathbf{x}) q(\mathbf{z})} d\mathbf{z} \\
&= \int q(\mathbf{z}) \ln \left(\frac{p(\boldsymbol{\theta}, \mathbf{z} | \mathbf{x})}{q(\mathbf{z})} \right) d\mathbf{z} + \int q(\mathbf{z}) \ln \frac{q(\mathbf{z})}{p(\mathbf{z} | \boldsymbol{\theta}; \mathbf{x})} d\mathbf{z}
\end{aligned}$$

Consequently this implies that EM algorithm still applies for the maximization of posterior with contingency on data and the corresponding $q(\mathbf{z})$ is replaced with $q(\mathbf{z} | \boldsymbol{\theta}^{\text{OLD}}; \mathbf{x})$, which essentially shares the same functional form with $q(\mathbf{z} | \mathbf{x}; \boldsymbol{\theta}^{\text{OLD}})$ at each iteration.

Remark 5.2 A direct implication from [Proposition 5.1](#) is that the objective $F(q, \boldsymbol{\theta})$ as the lower bound of log marginal likelihood $\ln p(\mathbf{x}; \boldsymbol{\theta})$, the gap between $\ln p(\mathbf{x}; \boldsymbol{\theta})$ and $F(q, \boldsymbol{\theta})$ is KL distance and hence maximizing ELBO minimizes KL, which is essentially the objective of Variational Bayes.

EM algorithm is essentially an iterative algorithm with parameters $\boldsymbol{\theta}$ constantly updated in a scheme such that $\boldsymbol{\theta}^{\text{OLD}} \rightarrow \boldsymbol{\theta}^{\text{NEW}} \rightarrow \boldsymbol{\theta}^{\text{OLD}} \dots$. At E-step, replacing $q(\mathbf{z})$ with $p(\mathbf{z} | \mathbf{x}; \boldsymbol{\theta}^{\text{OLD}})$ and substituting back to [\(21\)](#) yields

$$\begin{aligned}
F(\boldsymbol{\theta}, q) &= \int p(\mathbf{z} | \mathbf{x}; \boldsymbol{\theta}^{\text{OLD}}) \ln \left(\frac{p(\mathbf{x}, \mathbf{z} | \boldsymbol{\theta})}{p(\mathbf{z} | \mathbf{x}; \boldsymbol{\theta}^{\text{OLD}})} \right) d\mathbf{z} \\
&= \int p(\mathbf{z} | \mathbf{x}; \boldsymbol{\theta}^{\text{OLD}}) \ln (p(\mathbf{x}, \mathbf{z} | \boldsymbol{\theta})) d\mathbf{z} - \int p(\mathbf{z} | \mathbf{x}; \boldsymbol{\theta}^{\text{OLD}}) \ln (p(\mathbf{z} | \mathbf{x}; \boldsymbol{\theta}^{\text{OLD}})) d\mathbf{z} \\
&= Q(\boldsymbol{\theta}, \boldsymbol{\theta}^{\text{OLD}}) + \text{const}
\end{aligned}$$

Hence equivalently in the M-step, $\boldsymbol{\theta}$ is updated such that

$$\boldsymbol{\theta}^{\text{NEW}} = \arg \max_{\boldsymbol{\theta}} Q(\boldsymbol{\theta}, \boldsymbol{\theta}^{\text{OLD}}) \tag{22}$$

To sum up, EM algorithm essentially relies on the following two iterative steps repeatedly in a closed loop until difference between $\boldsymbol{\theta}^{\text{NEW}}$ and $\boldsymbol{\theta}^{\text{OLD}}$ is less than the tolerance rate specified.

E-step Compute $p(\mathbf{z} | \mathbf{x}; \boldsymbol{\theta}^{\text{OLD}})$ and $Q(\boldsymbol{\theta}, \boldsymbol{\theta}^{\text{OLD}})$

M-step Evaluate

$$\boldsymbol{\theta}^{\text{NEW}} = \arg \max_{\boldsymbol{\theta}} Q(\boldsymbol{\theta}, \boldsymbol{\theta}^{\text{OLD}})$$

Example 5.1 Consider a Gaussian mixture with G components indexed by $j = 1, \dots, G$ and for each component along with a specific observation \mathbf{x}_i we know exactly the density form $f_j(\mathbf{x}_i)$ and denote the parameters as the probabilities assigned to each component, π_j and $\sum_{j=1}^G \pi_j = 1$, collected in $\boldsymbol{\theta} = \{\pi_j\}_{j=1}^G$. All these lead to the log marginal likelihood of observed data \mathbf{x} as

$$\ln p(\mathbf{x} | \boldsymbol{\theta}) = \sum_{i=1}^N \ln \left[\sum_{j=1}^G \pi_j f_j(\mathbf{x}_i) \right] \quad (23)$$

An obvious difficulty for applying MLE with respect to this log marginal likelihood is the summation over g as the input of $\ln(\cdot)$. However this problem could be circumvented by introducing latent variable \mathbf{z} such that unobserved element contained in \mathbf{z} indexed by z_{ij} is 0-1 variable indicating whether the i -th unobserved component z_i associated with \mathbf{x}_i is from the j -th component. Consequently the joint log likelihood of \mathbf{x} and \mathbf{z} (sometime this is also referred to complete log likelihood) is

$$\ln p(\mathbf{x}, \mathbf{z} | \boldsymbol{\theta}) = \sum_{i=1}^N \sum_{j=1}^G z_{ij} [\ln \pi_j + \ln f_j(\mathbf{x}_i)] \quad (24)$$

Expectation of $\ln p(\mathbf{x}, \mathbf{z} | \boldsymbol{\theta})$ over \mathbf{z} taken with respect to $p(\mathbf{z} | \mathbf{x}; \boldsymbol{\theta}^{\text{OLD}})$, which is required for the **E-step**, is given as

$$\sum_{i=1}^N \sum_{j=1}^G \frac{\pi_j^{\text{OLD}} f_j(\mathbf{x}_i)}{\sum_{g=1}^G \pi_g^{\text{OLD}} f_g(\mathbf{x}_i)} [\ln \pi_j + \ln f_j(\mathbf{x}_i)] \quad (25)$$

Then for the **M-step**, we consider the following constrained optimization problem

$$\begin{aligned} \max \quad & \sum_{i=1}^N \sum_{j=1}^G \frac{\pi_j^{\text{OLD}} f_j(\mathbf{x}_i)}{\sum_{g=1}^G \pi_g^{\text{OLD}} f_g(\mathbf{x}_i)} [\ln \pi_j + \ln f_j(\mathbf{x}_i)] \\ \text{s. t.} \quad & \pi_1 + \dots + \pi_G = 1 \end{aligned}$$

the standard Lagrangian technique works here with corresponding multiplier denoted as λ and the related first order conditions for π_j implies that

$$\frac{1}{\pi_j} \sum_{i=1}^N \frac{\pi_j^{\text{OLD}} f_j(\mathbf{x}_i)}{\sum_{g=1}^G \pi_g^{\text{OLD}} f_g(\mathbf{x}_i)} - \lambda = 0, \quad 1 \leq j \leq G$$

and

$$\pi_j = \sum_{i=1}^N \frac{\pi_j^{OLD} f_j(\mathbf{x}_i)}{\sum_{g=1}^G \pi_g^{OLD} f_g(\mathbf{x}_i)} \Big/ \lambda$$

Taking summation over j on both sides of the above equation and using the required constraint solves λ as

$$\lambda = \sum_{i=1}^N \sum_{j=1}^G \frac{\pi_j^{OLD} f_j(\mathbf{x}_i)}{\sum_{g=1}^G \pi_g^{OLD} f_g(\mathbf{x}_i)} = N$$

Finally all the previous discussion yields that $\boldsymbol{\theta}$ updated from $\boldsymbol{\theta}^{OLD}$ to $\boldsymbol{\theta}^{NEW}$ in the **M-step** is as following

$$\boldsymbol{\theta}^{NEW} = \left(\sum_{i=1}^N \frac{\pi_1^{OLD} f_1(\mathbf{x}_i)}{\sum_{g=1}^G \pi_g^{OLD} f_g(\mathbf{x}_i)} \Big/ N, \dots, \sum_{i=1}^N \frac{\pi_G^{OLD} f_G(\mathbf{x}_i)}{\sum_{g=1}^G \pi_g^{OLD} f_g(\mathbf{x}_i)} \Big/ N \right) \quad (26)$$

A simple R code are provided for demonstrating the main logic of this procedure.

Example 5.2 Another direct application of EM algorithm is variable selection. Normally in Bayesian literature this target is achieved by imposing “spike-and-slab” prior (which is essentially a special case of Gaussian mixture) on coefficients and the corresponding importance weights of different variables are relied upon the posterior updates (George and McCulloch, 1993, 1997; Ročková and George, 2014). Specifically, we focus on the setting where there exists $n \times 1$ response vector (collection of observations of dependent variables), \mathbf{y} ; and potential variables (independent variables), $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_p]$, which is a $n \times p$ matrix. Temporarily a Gaussian linear model is assumed to relate \mathbf{X} to \mathbf{y} . That is, we generally have the following specified normal density,

$$f(\mathbf{y} | \mathbf{X}) = \mathcal{N}_n(\nu_n \boldsymbol{\alpha} + \mathbf{X}\boldsymbol{\beta}, \sigma^2 \mathbf{I}_n) \quad (27)$$

As with many Bayesian variable selection approaches, a binary latent variables $\boldsymbol{\gamma} = (\gamma_1, \dots, \gamma_p)^\top$ is introduced where $\gamma_i \in \{0, 1\}$ and $\gamma_i = 1$ indicates that x_i is included in model. The intrinsic logic for the variable selection objective is that in combination with some suitable prior for $\boldsymbol{\alpha}$, $\boldsymbol{\beta}$, σ and $\boldsymbol{\gamma}$, the induced posterior $\pi(\boldsymbol{\gamma} | \text{Data}) = \pi(\boldsymbol{\gamma} | \mathbf{y}, \mathbf{X})$ then summarizes all the postdata variable selection uncertainty. In particular, (George and McCulloch, 1997, henceforth GM97) proposed the following specified prior,

$$\pi(\boldsymbol{\beta} | \sigma, \boldsymbol{\gamma}, v_0, v_1) = \mathcal{N}_p(\mathbf{0}, D_{\sigma, \boldsymbol{\gamma}}) \quad (28)$$

where

$$D_{\sigma, \boldsymbol{\gamma}} = \sigma^2 \text{diag}(a_1, \dots, a_p).$$

with $a_i = (1 - \gamma_i)v_0 + \gamma_i v_1$ for $0 \leq v_0 < v_1$. Such a kind of specification implies that it is essentially a Gaussian mixture with two components. GM97 recommended setting v_0 and v_1 to be small and large values, respectively, to distinguish those β_i values which warrant exclusion of x_i from that inclusion of x_i . The key discussion corresponds closely to the specified distribution for $\pi(\boldsymbol{\gamma} | \boldsymbol{\theta})$. A natural

default choice for this would be the *i.i.d.* Bernoulli prior of the following form

$$\pi(\boldsymbol{\gamma} \mid \boldsymbol{\theta}) = \theta^{|\boldsymbol{\gamma}|} (1 - \theta)^{p-|\boldsymbol{\gamma}|} \quad (29)$$

where

$$\theta \in [0, 1] \text{ and } |\boldsymbol{\gamma}| = \sum_i \gamma_i.$$

For the prior of θ , normally we adopt beta distribution such that

$$\pi(\theta) \propto \theta^{a-1} (1 - \theta)^{b-1}$$

For the prior of σ^2 , [GM97](#) proposes the use of inverse gamma prior

$$\pi(\sigma^2 \mid \boldsymbol{\gamma}) = \text{IG}\left(\frac{\nu}{2}, \frac{\nu\lambda}{2}\right)$$

where this inverse gamma prior takes the following functional PDF form ³,

$$\frac{(\nu\lambda/2)^{\nu/2}}{\Gamma(\nu/2)} (\sigma^2)^{-\nu/2-1} \exp\left(-\frac{\nu\lambda}{2\sigma^2}\right)$$

Note that for any practical application of EM algorithm, our target is never to exactly pin down the posterior distribution of latent variables but the posterior updates of parameters instead. Consequently given our setting here, parameters are collected in $(\boldsymbol{\beta}, \boldsymbol{\theta}, \sigma)$ ⁴ and our ultimate target is geared toward finding **posterior modes of parameter posterior** $\pi(\boldsymbol{\beta}, \boldsymbol{\theta}, \sigma \mid \mathbf{y}, \mathbf{X})$ rather than simulating from the entire **model posterior** $\pi(\boldsymbol{\gamma} \mid \mathbf{y}, \mathbf{X})$.

Following the conventional procedure implemented in EM algorithm, maximization of $\pi(\boldsymbol{\beta}, \boldsymbol{\theta}, \sigma \mid \mathbf{y}, \mathbf{X})$ is implemented indirectly by proceeding iteratively in terms of the “complete log likelihood” as we did in the previous example, that is, $\ln \pi(\boldsymbol{\beta}, \boldsymbol{\theta}, \sigma, \boldsymbol{\gamma} \mid \mathbf{y}, \mathbf{X})$, where the inclusion indicators collected in $\boldsymbol{\gamma}$ are treated as “missing data”. More precisely, the whole procedure comprises the iterative maximization of the following objective function

$$\begin{aligned} & \mathbf{Q}(\boldsymbol{\beta}, \boldsymbol{\theta}, \sigma \mid \boldsymbol{\beta}^{(k)}, \boldsymbol{\theta}^{(k)}, \sigma^{(k)}) \\ &= \mathbb{E}_{\boldsymbol{\gamma} \mid \cdot} \left[\ln \pi(\boldsymbol{\beta}, \boldsymbol{\theta}, \sigma, \boldsymbol{\gamma} \mid \mathbf{y}, \mathbf{X}) \mid \boldsymbol{\beta}^{(k)}, \boldsymbol{\theta}^{(k)}, \sigma^{(k)}, \mathbf{y} \right] \end{aligned} \quad (30)$$

where the operator $\mathbb{E}_{\boldsymbol{\gamma} \mid \cdot}(\cdot)$ denotes the conditional expectation, that is,

$$\mathbb{E}_{\boldsymbol{\gamma} \mid \boldsymbol{\beta}^{(k)}, \boldsymbol{\theta}^{(k)}, \sigma^{(k)}, \mathbf{y}, \mathbf{X}}(\cdot)$$

³ In general for inverse-gamma distribution with parameter α and β , the corresponding PDF (probability density function) takes the form of $\frac{\beta^\alpha}{\Gamma(\alpha)} x^{-\alpha-1} \exp(-\frac{\beta}{x})$, where x refers to the random variable.

⁴ Here we ignore the emphasis of α since generally it can be subsumed in $\boldsymbol{\beta}$ by treating it as the coefficient of fixed constant.

Specifically, (30) can be decomposed into the addition of a constant term and the following additional two terms⁵

$$\begin{aligned} & \mathbf{Q}_1(\boldsymbol{\beta}, \sigma \mid \boldsymbol{\beta}^{(k)}, \boldsymbol{\theta}^{(k)}, \sigma^{(k)}) \\ &= -\frac{(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^\top (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})}{2\sigma^2} - \frac{n+p+v+2}{2} \ln(\sigma^2) - \frac{\nu\lambda}{2\sigma^2} - \frac{1}{2\sigma^2} \sum_{i=1}^p \beta_i^2 \mathbb{E}_{\boldsymbol{\gamma} \mid \cdot} \left[\frac{1}{v_0(1-\gamma_i) + v_1\gamma_i} \right] \end{aligned}$$

$$\begin{aligned} & \mathbf{Q}_2(\boldsymbol{\theta} \mid \boldsymbol{\beta}^{(k)}, \boldsymbol{\theta}^{(k)}, \sigma^{(k)}) \\ &= \sum_{i=1}^p \ln\left(\frac{\theta}{1-\theta}\right) \mathbb{E}_{\boldsymbol{\gamma} \mid \cdot} \gamma_i + (a-1) \ln(\theta) + (p+b-1) \ln(1-\theta) \end{aligned}$$

Thus in general we have

$$\mathbf{Q}(\boldsymbol{\beta}, \boldsymbol{\theta}, \sigma \mid \boldsymbol{\beta}^{(k)}, \boldsymbol{\theta}^{(k)}, \sigma^{(k)}) = \mathbf{C} + \mathbf{Q}_1(\boldsymbol{\beta}, \sigma \mid \boldsymbol{\beta}^{(k)}, \boldsymbol{\theta}^{(k)}, \sigma^{(k)}) + \mathbf{Q}_2(\boldsymbol{\theta} \mid \boldsymbol{\beta}^{(k)}, \boldsymbol{\theta}^{(k)}, \sigma^{(k)}) \quad (31)$$

Such a kind of separability into two distinct functions yields a M -step that is obtained by optimizing each function. To see why $\mathbf{Q}_1(\cdot)$ and $\mathbf{Q}_2(\cdot)$ take the demonstrated functional form respectively, firstly we observe the generic fact that any well-defined posterior is proportional to the corresponding jointly likelihood up to a constant which is the result of marginalizing out the parameters. The everything remained is just obtained by taking log of the product of prior and the likelihood of data given specified parameter. Analogous discussion applies for $\mathbf{Q}_2(\cdot)$. More specifically, for $\mathbf{Q}_1(\cdot)$, we have

$$\begin{array}{ccc} f(\mathbf{y} \mid \mathbf{X}) \cdot \pi(\boldsymbol{\beta} \mid \sigma, \boldsymbol{\gamma}, v_0, v_1) \cdot \pi(\sigma^2 \mid \boldsymbol{\gamma}) \\ \uparrow \qquad \qquad \uparrow \qquad \qquad \uparrow \\ \mathcal{N}_n(\nu_n \boldsymbol{\alpha} + \mathbf{X}\boldsymbol{\beta}, \sigma^2 \mathbf{I}_n) \quad \mathcal{N}_p(\mathbf{0}, D_{\sigma, \boldsymbol{\gamma}}) \quad \text{IG}\left(\frac{\nu}{2}, \frac{\nu\lambda}{2}\right) \end{array}$$

Taking log of the above product yields the desired functional form of $\mathbf{Q}_1(\cdot)$; For $\mathbf{Q}_2(\cdot)$, we have

$$\begin{aligned} & \pi(\boldsymbol{\theta}) \cdot \pi(\boldsymbol{\gamma} \mid \boldsymbol{\theta}) \\ & \propto \theta^{a-1} (1-\theta)^{b-1} \theta^{|\boldsymbol{\gamma}|} (1-\theta)^{p-|\boldsymbol{\gamma}|} = \theta^{a-1} (1-\theta)^{p+b-1} \prod_{i=1}^p \left(\frac{\theta}{1-\theta}\right)^{\gamma_i} \end{aligned}$$

Taking log of the above formula yields the desired functional form of $\mathbf{Q}_2(\cdot)$

$$(a-1) \ln(\theta) + (p+b-1) \ln(1-\theta) + \sum_{i=1}^p \ln\left(\frac{\theta}{1-\theta}\right) \gamma_i$$

hence in the E -step where γ_i is replaced with its corresponding conditional expectation, $\mathbb{E}_{\boldsymbol{\gamma} \mid \cdot} \gamma_i$, we

⁵ It is a little bit different from the formula derived in GM97, where the multiplier before $\ln(\sigma^2)$ is given as $\frac{n-1+p+\nu}{2}$. I personally hold the opinion that the formula given in GM97 should be a typo.

obtain the desired functional form of $\mathbf{Q}_2(\cdot)$. From the previous discussion, two conditional expectation plays the vital role in proceeding the E-step, thus $\mathbb{E}_{\gamma} \mid \cdot \left[\frac{1}{v_0(1-\gamma_i)+v_1\gamma_i} \right]$ and $\mathbb{E}_{\gamma} \mid \cdot \gamma_i$ respectively. For this reason we briefly discuss how this two conditional expectation is calculated,

- $\mathbb{E}_{\gamma} \mid \cdot \gamma_i$:

$$\mathbb{E}_{\gamma} \mid \cdot \gamma_i = \mathbf{P} \left(\gamma_i = 1 \mid \boldsymbol{\beta}^{(k)}, \boldsymbol{\theta}^{(k)}, \sigma^{(k)} \right) = p_i^*$$

where

$$p_i^* = \frac{a_i}{a_i + b_i}$$

and

$$\begin{aligned} a_i &= \pi \left(\beta_i^{(k)} \mid \sigma^{(k)}, \gamma_i = 1 \right) \mathbf{P} \left(\gamma_i = 1 \mid \boldsymbol{\theta}^{(k)} \right) \\ b_i &= \pi \left(\beta_i^{(k)} \mid \sigma^{(k)}, \gamma_i = 0 \right) \mathbf{P} \left(\gamma_i = 0 \mid \boldsymbol{\theta}^{(k)} \right) \end{aligned}$$

- $\mathbb{E}_{\gamma} \mid \cdot \left[\frac{1}{v_0(1-\gamma_i)+v_1\gamma_i} \right]$:

$$\begin{aligned} &\mathbb{E}_{\gamma} \mid \cdot \left[\frac{1}{v_0(1-\gamma_i)+v_1\gamma_i} \right] \\ &= \frac{\mathbb{E}_{\gamma} \mid \cdot (1-\gamma_i)}{v_0} + \frac{\mathbb{E}_{\gamma} \mid \cdot \gamma_i}{v_1} = \frac{1-p_i^*}{v_0} + \frac{p_i^*}{v_1} \equiv d_i^* \end{aligned}$$

One attractive feature of the previously discussed setting is that all the required maximization in M-step yields analytical solution so that $\boldsymbol{\beta}^{(k)}$, $\boldsymbol{\theta}^{(k)}$ and $\sigma^{(k)}$ could be updated quickly to $\boldsymbol{\beta}^{(k+1)}$, $\boldsymbol{\theta}^{(k+1)}$ and $\sigma^{(k+1)}$ with $\boldsymbol{\beta}^{(k+1)}$, $\boldsymbol{\theta}^{(k+1)}$ and $\sigma^{(k+1)}$ replaced by the corresponding optimal solution respectively. Under some regular conditions, this sequence would converge at $\widehat{\boldsymbol{\beta}}$, $\widehat{\boldsymbol{\theta}}$ and $\widehat{\sigma}$. Then based on these estimated parameters, the conditional component inclusion probabilities are given by

$$\mathbf{P} \left(\gamma_i \mid \widehat{\beta}_i, \widehat{\boldsymbol{\theta}}, \widehat{\sigma} \right) = \frac{c_i}{c_i + d_i} \quad (32)$$

where

$$\begin{aligned} c_i &= \pi \left(\widehat{\beta}_i, \widehat{\sigma}, \gamma_i = 1 \right) \pi \left(\gamma_i = 1 \mid \widehat{\boldsymbol{\theta}} \right) \\ d_i &= \pi \left(\widehat{\beta}_i, \widehat{\sigma}, \gamma_i = 0 \right) \pi \left(\gamma_i = 0 \mid \widehat{\boldsymbol{\theta}} \right) \end{aligned}$$

Definition 2 α -Rényi divergence is defined for measuring the relative divergence between two probability measures. Suppose $\alpha \in (0, 1)$ and P and Q are two probability measures. Let μ be any measure such that $P \ll \mu$ and $Q \ll \mu$, for example $\mu = P + Q$, then α -Rényi divergence is defined as following

$$D_{\alpha}(P, Q) = \frac{1}{\alpha - 1} \log \int \left(\frac{dP}{d\mu} \right)^{\alpha} \left(\frac{dQ}{d\mu} \right)^{1-\alpha} d\mu \quad (33)$$

Comprehensive discussion about α -Rényi divergence is available in Van Erven and Harremoës (2014). An important convergence would be as following,

$$\begin{aligned}\lim_{\alpha \rightarrow 0} \frac{1}{\alpha - 1} \ln \int p^\alpha q^{1-\alpha} d\mu &= -\ln \int \mathbf{1}_{\{p>0\}} q d\mu = -\ln Q (p > 0) \\ \lim_{\alpha \rightarrow 1} \frac{1}{\alpha - 1} \ln \int p^\alpha q^{1-\alpha} d\mu &= \int p \ln \left(\frac{p}{q} \right) d\mu\end{aligned}$$

where $p = \frac{dP}{d\mu}$ and $q = \frac{dR}{d\mu}$.

Remark 5.3 Essentially the justification for these two limits is based on the justification for commuting the integration and the limit taken with respect to α . This is carefully discussed in Van Erven and Harremoës (2014), where **Lemma 1** and the corresponding derivation for **Theorem 5** jointly justify such a kind of desired communication. The fundamental tool used for proving these claims is monotone convergence theorem. Rather than discussing the proof rigorously here, an important and useful observation would be that $\frac{d}{d\alpha} p^\alpha q^{1-\alpha} = p^\alpha q^{1-\alpha} \ln \left(\frac{p}{q} \right)$, accordingly once communication between integration and limit taken with respect to α is rigorously justified, a heuristic demonstration of why the claimed limits hold is by applying L'Hôpital's rule such that

$$\lim_{\alpha \rightarrow 1} \frac{1}{\alpha - 1} p^\alpha q^{1-\alpha} = \lim_{\alpha \rightarrow 1} p^\alpha q^{1-\alpha} \ln \left(\frac{p}{q} \right) = p \ln \left(\frac{p}{q} \right).$$

Another useful facts are as long as $P \ll Q$, then we have

- $p < q$ almost surely holds with respect to Q .
- If we define $x_\alpha = \int p^\alpha q^{1-\alpha} d\mu$, then it is appropriate to have $\lim_{\alpha \rightarrow 1} \int p^\alpha q^{1-\alpha} d\mu = \int p d\mu = 1$. More specifically, this convergence is depicted from different directions as following

$$\int p^\alpha q^{1-\alpha} d\mu = \int \left(\frac{p}{q} \right)^\alpha q d\mu$$

$$\alpha \uparrow 1 \quad x_\alpha \downarrow 1 \quad x_\alpha - 1 \geq \ln x_\alpha \rightarrow 0$$

$$\alpha \downarrow 1 \quad x_\alpha \uparrow 1 \quad \ln x_\alpha \geq x_\alpha - 1 \rightarrow 0$$

where $\uparrow 1$ refers to the convergence below 1 while $\downarrow 1$ refers to the convergence above 1.

Consequently,

$$\lim_{\alpha \rightarrow 1} \frac{1}{\alpha - 1} \ln \int p^\alpha q^{1-\alpha} d\mu = \lim_{\alpha \rightarrow 1} \frac{\int p^\alpha q^{1-\alpha} d\mu - 1}{\alpha - 1} = \lim_{\alpha \rightarrow 1} \int_{p,q>0} \frac{p - p^\alpha q^{1-\alpha}}{1 - \alpha} d\mu \quad (34)$$

and this implies as long as we are able to change the order of limits and integration and then taking limit with respect to α by applying L'Hôpital's rule yields

$$\lim_{\alpha \rightarrow 1} D_\alpha (P, Q) = \text{KL} (P, Q) \quad (35)$$

With this introduced α -Rényi divergence, we introduce the following concepts for later discussion

- Let Q denote the dominating measure for the family of distributions for data and hence the associated Radon-Nikodym derivative (probability density function in some sense) is as following

$$p_{\boldsymbol{\theta}} = \frac{dP_{\boldsymbol{\theta}}}{dQ}$$

- $\boldsymbol{\theta} \in \Theta$ and Θ is equipped with proper σ -algebra \mathcal{T} , $\mathcal{M}_1^+(\Theta, \mathcal{T})$ refers to the set of all probability distributions on measurable space (Θ, \mathcal{T}) . Accordingly, $\pi \in \mathcal{M}_1^+(\Theta, \mathcal{T})$ denotes prior.
- Likelihood as the function of $\boldsymbol{\theta}$ is ⁶

$$L_n(\boldsymbol{\theta}) = \prod_{i=1}^n p_{\boldsymbol{\theta}}(x_i)$$

- $\forall(\boldsymbol{\theta}, \boldsymbol{\theta}')$, negative log-likelihood ratio is defined as following

$$r_n(\boldsymbol{\theta}, \boldsymbol{\theta}') = \sum_{i=1}^n \log \frac{p_{\boldsymbol{\theta}'}(x_i)}{p_{\boldsymbol{\theta}}(x_i)}$$

- Fractional posterior is then defined as following

$$\pi_{n,\alpha}(d\boldsymbol{\theta} | \mathbf{x}_1^n) = \frac{e^{-\alpha r_n(\boldsymbol{\theta}, \boldsymbol{\theta}_0)} \pi(d\boldsymbol{\theta})}{\int e^{-\alpha r_n(\boldsymbol{\theta}, \boldsymbol{\theta}_0)} \pi(d\boldsymbol{\theta}) \pi(d\boldsymbol{\theta})} \quad (36)$$

With these introduced concepts and notations, we summarize some important results discussed in Alquier and Ridgway (2020):

Theorem 5.1 (PAC-Bayesian Inequality) For any $\alpha \in (0, 1)$ and any $\epsilon \in (0, 1)$,

$$\begin{aligned} \mathbb{P} \left(\forall \rho \in \mathcal{M}_1^+(\Theta), \int D_{\alpha}(P_{\boldsymbol{\theta}}, P_{\boldsymbol{\theta}_0}) \rho(d\boldsymbol{\theta}) \right. \\ \left. \leq \frac{\alpha}{1-\alpha} \int \frac{r_n(\boldsymbol{\theta}, \boldsymbol{\theta}_0)}{n} \rho(d\boldsymbol{\theta}) + \frac{\text{KL}(\rho, \pi) + \log(1/\epsilon)}{n(1-\alpha)} \right) \geq 1 - \epsilon \end{aligned} \quad (37)$$

This result originates from Catoni (2004, 2007) and is extended to variational approximation by alternatively defining the following approximate posterior as

$$\begin{aligned} \tilde{\pi}_{n,\alpha}(\cdot | \mathbf{x}_1^n) &= \operatorname{argmin}_{\rho \in \mathcal{F}} \left\{ \alpha \int r_n(\boldsymbol{\theta}, \boldsymbol{\theta}_0) \rho(d\boldsymbol{\theta}) + \text{KL}(\rho, \pi) \right\} \\ &= \operatorname{argmin}_{\rho \in \mathcal{F}} \left\{ -\alpha \int \sum_{i=1}^n \log p_{\boldsymbol{\theta}}(x_i) \rho(d\boldsymbol{\theta}) + \text{KL}(\rho, \pi) \right\} \end{aligned}$$

⁶ x_i refers to the i -th realization of random variables (either scalars or vectors) \mathbf{x} .

Corollary 5.1 (Concentration of VB approximation) For any $\alpha \in (0, 1)$ and $\epsilon \in (0, 1)$, with probability at least $1 - \epsilon$

$$\int D_\alpha(P_\theta, P_{\theta_0}) \tilde{\pi}_{n,\alpha}(\mathrm{d}\theta \mid \mathbf{x}_1^n) \leq \inf_{\rho \in \mathcal{F}} \left\{ \frac{\alpha}{1-\alpha} \int \frac{r_n(\theta, \theta_0)}{n} \rho(\mathrm{d}\theta) + \frac{\mathrm{KL}(\rho, \pi) + \log(1/\epsilon)}{n(1-\alpha)} \right\} \quad (38)$$

[Theorem 5.1](#) and [Corollary 5.1](#) jointly leads to the following claimed main result in Alquier and Ridgway (2020)

Theorem 5.2 For any fixed $\mathcal{F} \subset \mathcal{M}_1^+(\Theta)$, if we assume that there exists $\varepsilon_n > 0$ such that

$$\int \mathrm{KL}(P_\theta, P_{\theta_0}) \rho_n(\mathrm{d}\theta) \leq \varepsilon_n, \quad \int \mathbb{E} \left[\log^2 \left(\frac{p_\theta(\mathbf{x}_i)}{p_{\theta_0}(\mathbf{x}_i)} \right) \right] \rho_n(\mathrm{d}\theta) \leq \varepsilon_n$$

and

$$\mathrm{KL}(\rho_n, \pi) \leq n\varepsilon_n$$

Then for any $\alpha \in (0, 1)$, for any $(\varepsilon, \eta) \in (0, 1)^2$, the following inequality holds

$$\mathbb{P} \left[\int D_\alpha(P_{\theta_0}, P_\theta) \tilde{\pi}_{n,\alpha}(\mathrm{d}\theta \mid \mathbf{x}_1^n) \leq \frac{(\alpha+1)\varepsilon + \alpha \sqrt{\frac{\varepsilon_n}{n\eta}} + \frac{\log(\frac{1}{\varepsilon})}{n}}{1-\alpha} \right] \geq 1 - \varepsilon - \eta \quad (39)$$

Remark 5.4 There are some remarks to be emphasized here for discussing this main result.

- How ρ_n is selected. Define the following restricted area

$$B(r) = \left\{ \theta \in \Theta : \mathrm{KL}(P_{\theta_0}, P_\theta) \leq r, \mathbb{E} \left[\log^2 \left(\frac{p_\theta(\mathbf{x}_i)}{p_{\theta_0}(\mathbf{x}_i)} \right) \right] \leq r \right\}$$

and then ρ_n is selected such that

$$\rho_n = \pi|_{B(\varepsilon_n)}$$

i.e. ρ_n is selected as the π restricted to $B(\varepsilon_n)$ and hence for this case the required condition rewrites as following

$$\mathrm{KL}(\rho_n, \pi) = -\log \pi(B(\varepsilon_n)) \leq n\varepsilon_n$$

- An alternative explanation for the results implied from this theorem is that Expected (with respect to VB posterior) α -Rényi divergence of data distribution associated with different θ cannot be that divergent with probability $1 - \varepsilon - \eta$ and this probability is with respect to data.

Alternatively as we have mentioned earlier that we are allowed to regard θ as latent variables but temporarily treat it as the fixed parameter. By doing so, it is appropriate to redefine the evidence lower bound (ELBO) by simultaneously accounting for both local latent variables \mathbf{z} and

global latent variables $\boldsymbol{\theta}$. This trick treating the global variable $\boldsymbol{\theta}$ as temporarily fixed also justifies our discussion temporarily within frequentist domain and later the following discussion about how Variational Frequentist and Variational Bayes are connected. To formally set up the connection, it would be helpful to summarize the associated definitions as following

[Place Table 1 about here]

5.2 Variational Expectation Maximization

The key difference between Variational Expectation Maximization (VEM) and Expectation Maximization (EM) is mainly at E-step where there is an implicitly crucial assumption that $p(\mathbf{z} | \mathbf{x}; \boldsymbol{\theta}^{\text{OLD}})$ is available so that $q(\mathbf{z})$ can be replaced with $p(\mathbf{z} | \mathbf{x}; \boldsymbol{\theta}^{\text{OLD}})$. However this is not generally the case when $p(\mathbf{z} | \mathbf{x}; \boldsymbol{\theta}^{\text{OLD}})$ is tractable, hence we turn to find the optimal functionality of $q(\cdot)$ instead by relying on the assumption that $q(\mathbf{z}) = \prod_i q(z_i)$, then ELBO can be factorized as following

$$\begin{aligned}
F(q, \boldsymbol{\theta}) &= \int q(\mathbf{z}) \ln \left(\frac{p(\mathbf{z}, \mathbf{x}; \boldsymbol{\theta})}{q(\mathbf{z})} \right) d\mathbf{z} \\
&= \int q(z_j) \left(\int \prod_{i \neq j} q(z_i) \ln p(\mathbf{z}, \mathbf{x}; \boldsymbol{\theta}) \right) \prod_{i \neq j} dz_i dz_j - \int q(z_j) \ln q(z_j) dz_j - \sum_{i \neq j} \int q(z_i) \ln q(z_i) dz_i \\
&= \int q(z_j) \ln \left(\frac{\exp(\langle \ln p(\mathbf{z}, \mathbf{x} | \boldsymbol{\theta}) \rangle_{i \neq j})}{q(z_j)} \right) dz_j - \sum_{i \neq j} \int q(z_i) \ln q(z_i) dz_i \\
&= \int q(z_j) \ln \left(\frac{\tilde{p}_{i \neq j}}{q(z_j)} \right) dz_j + H(z_{i \neq j}) + \text{const} \\
&= -\text{KL}(q_j \| \tilde{p}_{i \neq j}) + H(z_{i \neq j}) + \text{const}
\end{aligned}$$

where $\tilde{p}_{i \neq j}$ is normalized pdf. Since K-L divergence is non-negative, with $z_{i \neq j}$ to be integrated out under proper expectation measure, the sufficient and necessary condition for ELBO to be maximized is $\text{KL}(q_j \| \tilde{p}_{i \neq j}) = 0$, which sufficiently holds when

$$q(z_j) = \tilde{p}_{i \neq j} \propto \exp(\langle \ln p(\mathbf{z}, \mathbf{x} | \boldsymbol{\theta}) \rangle_{i \neq j})$$

or in other words the log of optimal density $q_j^*(z_j)$ is

$$\ln q_j^*(z_j) = \langle \ln p(\mathbf{x}, \mathbf{z}; \boldsymbol{\theta}) \rangle_{i \neq j} + \text{const} \quad (40)$$

Consequently, the EM algorithm is given by the following two steps

Variational E-step Evaluate $q^{\text{NEW}}(\mathbf{z})$ to maximize $F(q, \boldsymbol{\theta}^{\text{OLD}})$ solving (40)

Variational M-step Find

$$\theta^{\text{NEW}} = \arg \max_{\theta} F(\mathbf{q}^{\text{NEW}}, \theta)$$

In comparison to conventional data used for economics and finance research, information encoded in textual data is relatively more rich and serves as good complement to the conventional structured data. Recently there seems to be explosion of social science researches using text data directly or text data extracted from videos. One reason that there were relatively less empirical researches in this norm is the intrinsic features that data of these kinds are high dimensional while most of the analysis tool is not as powerful as the available analysis tools nowadays. Conventionally, classical MCMC methods serve as the major tools for the application of Bayesian methods in practice but one fatal problems encountered in the MCMC is its deteriorated computational efficiency especially for the case when both the covariates dimension and sample size are large. Variational Bayes (VB) by comparison replaces the conventional MCMC sampling procedure with optimization problem described previously so that to some extent alleviates the computational inefficiency, and VB is nowadays is increasingly becoming the major tools for many computationally demanding applications such as image and video processing, and natural language processing (NLP).

Unstructured data recorded in texts, videos can reveal some useful information for financial market and one important mechanism through which the unstructured data play the pivotal role is that the people’s emotion revealed from data is contagious and hence generates impact on people’s decisions corresponding to economic activities, this is also the mechanism referred by some existing literature as social transmission. To some extent, such a hypothesized mechanism is quite natural since psychologically people’s subjective perceptions are easily influenced by others so are associated economic behaviours.

Development in modern information technologies and computer science provides analysis tools both convenient and powerful to extract information collected in aforementioned unstructured data like texts and videos. Platforms like Google Cloud and Microsoft Azure both provide extensive support for modern data science language such as R and Python for natural language processing. Besides, community within the ecosystems of both R and Python have contributed a lot of integrated convenient packages such as `tm`, `topicmodels`, `googlelanguageR`, `tuneR`, `av` for R and `youtube_dl` for Python. Gentzkow, Kelly, and Tady (2019) recently provides a relatively comprehensive review about how the recently developed techniques have been applied for analysing textual data in economics and finance.

As what is mentioned previously, data extracted from text or video is highly unstructured and hence to make it manageable, some preliminary procedures have to be implemented. One important measure commonly used in textual analysis for this purpose is called “term frequency-inverse document frequency” (abbreviated as `tf-idf`). Generally, it is defined as following

$$\text{tf-idf}: \text{tf}_{ij} \times \text{idf}_j \tag{41}$$

where \mathbf{tf}_{ij} refers to the term frequency of word j in document i such that

$$\mathbf{tf}_{ij} = \frac{c_{ij}}{\sum_k c_{ik}}$$

with c_{ij} defined as the count of occurrences of word j in document i ; And \mathbf{idf}_j refers to: $\log(n/d_j)$ with

n : total number of documents

d_j : $\sum_i \mathbf{1}(c_{ij} > 0)$, $\mathbf{1}(\cdot)$ is indicator function.

Since in general we want to make data extracted from text reveal distinguishable information, those words rarely appearing in each document and commonly appearing in most or all of the documents should be eliminated. This can be achieved by setting a threshold value for the $\mathbf{tf-idf}$ measure since by construction small \mathbf{tf}_{ij} implies that the j -th word is rare in document i and similarly small \mathbf{idf}_j implies that the j -th word is relatively much common in all the documents. Essentially $\mathbf{tf-idf}$ as the transformed measure defines the most representative terms in a given document to be those that appear infrequently overall, but frequently in that specific document (see Engle et al., 2020, henceforth EGLKS2020). Based on the $\mathbf{tf-idf}$ measure, EGLKS2020 constructs a climate change index by calculating the “cosine similarity” between $\mathbf{tf-idf}$ measures for “climate change vocabulary (CCV)” and each daily *The Wall Street Journal* (WSJ) edition. Another index referred to Crimson Hexagon’s negative sentiment climate change index is constructed along with climate change index as well. One reason motivating the construction of Crimson Hexagon’s (CH) negative sentiment climate change index in companion with climate change index for comparison, as articulated in EGLKS2020, is that the construction of WSJ Climate Change News Index implicitly embeds the assumption that there is no discrepancy between good news and bad news about climate changes given the coverage of all WSJ corpus, which is potentially at the risk of inaccurately capturing the positive news about climate change. Consequently, Crimson Hexagon’s (CH) negative sentiment climate change index as one alternative for addressing this concern is forwarded in EGLKS2020 by confining the focus specifically on negative climate change news. Particularly, the CH Negative Climate Change Index is constructed as the “share of all news articles that are both about ‘climate change’ and that have been assigned to the ‘negative sentiment’ category.”

5.2.1 Discussion about parametrized optimization

It is obvious from the previous discussion that optimization plays the pivotal role in the implementation of VB method and consequently we may briefly discuss the general optimization procedure we often implement in practice. Stochastic gradient ascent method proposed in Robbins and Monro (1951) is widely used in practice for solving the optimization in **M-step**. For notation simplicity and consistency while emphasizing the associated dependency of ELBO on parameters collected in

ϕ , we denote the parametrized ELBO as

$$\mathcal{F}(\phi) = F(q, \phi, \theta) = \int q_\phi(\mathbf{z}) \ln \left(\frac{p(\mathbf{x}, \mathbf{z} | \theta)}{q_\phi(\mathbf{z})} \right) d\mathbf{z}$$

where we emphasize that $\mathcal{F}(\cdot)$ only depends on the ϕ with θ fixed at constant and then we are able to apply stochastic approximation method by calculating the unbiased gradient estimates $\widehat{\nabla_\phi \mathcal{F}(\phi)}$. To obtain the explicit formula for ∇_ϕ , it would be useful to note a key fact that

$$\begin{aligned} \mathbb{E}_q [\nabla_\phi \ln q_\phi(\mathbf{z})] &= \mathbb{E}_q \left[\frac{1}{q_\phi(\mathbf{z})} \nabla_\phi q_\phi(\mathbf{z}) \right] \\ &= \int q_\phi(\mathbf{z}) \frac{1}{q_\phi(\mathbf{z})} \nabla_\phi q_\phi(\mathbf{z}) d\mathbf{z} \\ &= \int \nabla_\phi q_\phi(\mathbf{z}) d\mathbf{z} = 0 \quad \left(\text{since } \int q_\phi(\mathbf{z}) d\mathbf{z} = 1 \right) \end{aligned}$$

This fact is also referred to the so-called log-derivative trick. With this observed fact, we are able to derive as following

$$\begin{aligned} \nabla_\phi \mathcal{F}(\phi) &= \int \nabla_\phi q_\phi(\mathbf{z}) [\ln p(\mathbf{x}, \mathbf{z} | \theta) - \ln q_\phi(\mathbf{z})] d\mathbf{z} - \int q_\phi(\mathbf{z}) \frac{1}{q_\phi(\mathbf{z})} \nabla_\phi q_\phi(\mathbf{z}) d\mathbf{z} \\ &= \mathbb{E}_q [\nabla_\phi \ln q_\phi(\mathbf{z}) (\ln p(\mathbf{x}, \mathbf{z} | \theta) - \ln q_\phi(\mathbf{z}))] \end{aligned} \tag{42}$$

since gradient calculated as in (42) is expectation taken with respect to $q_\phi(\cdot)$, it is easy to calculate the estimated gradient $\widehat{\nabla_\phi \mathcal{F}(\phi)}$ unbiasedly using samples from $q_\phi(\mathbf{z})$ provided the analytical form of $q_\phi(\mathbf{z})$ and the corresponding sampling from $q_\phi(\mathbf{z})$ is readily available. Then the basic idea of stochastic gradient descent method is as following. With $\mathcal{F}(\phi)$ as objective function to optimized and assume that under some regular conditions $\mathcal{F}(\phi)$ is differentiable, thus $\nabla_\phi \mathcal{F}(\phi)$ is its gradient and $\widehat{\nabla_\phi \mathcal{F}(\phi)}$, we start from an initial value $\phi^{(0)}$ and implement the following recursion for $t = 0, 1, \dots$

$$\phi^{(t+1)} = \phi^{(t)} + \rho_t \circ \widehat{\nabla_\phi \mathcal{F}(\phi^{(t)})}$$

until a certain stopping criteria is satisfied. It can be justified that as long as $\rho_t, t \geq 0$ (learning rate) satisfies Robbins-Monro conditions, that is $\sum_t \rho_t = \infty$ and $\sum_t \rho_t^2 < \infty$, then convergence of the sequence $\phi^{(t)}$ will be to local optimum. There are many discussions on choosing learning rate ρ_t in literature and ADADELTA method (Zeiler, 2012) is one relatively easy to implement and hence widely used in practice. Specifically for this setting at iteration $t + 1$, the i -th element ϕ_i of ϕ is updated as

$$\phi_i^{(t+1)} = \phi_i^{(t)} + \Delta \phi_i^{(t)}$$

where the incremental step size $\Delta\phi_i^{(t)}$ is $\rho_i^{(t)} g_{\phi_i}^{(t)}$ and $g_{\phi_i}^{(t)}$ denotes the i -th component $\widehat{\nabla_{\phi} \mathcal{F}(\phi^{(t)})}$ with $\rho_i^{(t)} g_{\phi_i}^{(t)}$ adaptively selected as

$$\rho_i^{(t)} = \frac{\sqrt{\mathbb{E}(\Delta_{\phi_i}^2)^{(t-1)} + \epsilon}}{\sqrt{\mathbb{E}(g_{\phi_i}^2)^{(t)} + \epsilon}}$$

where ϵ is a small positive constant and $\mathbb{E}(\Delta_{\phi_i}^2)^{(t)}$, $\mathbb{E}(g_{\phi_i}^2)^{(t)}$ are decayed running average of $\Delta\phi_i^{(t)2}$ and $g_{\phi_i}^{(t)2}$ respectively, defined by

$$\begin{aligned}\mathbb{E}(\Delta_{\phi_i}^2)^{(t)} &= \zeta \mathbb{E}(g_{\phi_i}^2)^{(t-1)} + (1 - \zeta) \Delta\phi_i^{(t)2} \\ \mathbb{E}(g_{\phi_i}^2)^{(t)} &= \zeta \mathbb{E}(g_{\phi_i}^2)^{(t-1)} + (1 - \zeta) g_{\phi_i}^{(t)2}\end{aligned}$$

Typically tuning parameters are specified as $\epsilon = 10^{-6}$ and $\zeta = 0.95$ the recursion is initialized at $\mathbb{E}(\Delta_{\phi_i}^2)^{(t)} = \mathbb{E}(g_{\phi_i}^2)^{(t)} = 0$.

Another way for implementing optimization is by directly minimizing the divergence between the approximate distribution and target posterior distribution. For discussion simplicity, introducing the following shorthand notation

$$\mathcal{D}(q) := \mathcal{D}(q, p_{\mathbf{x}}) \tag{43}$$

$$\tilde{\mathcal{D}}(q) := \mathcal{D}(q, f) \tag{44}$$

where $p_{\mathbf{x}}$ as the shorthand notation is in consistence with the aforementioned target posterior $p(\mathbf{z} | \mathbf{x}; \boldsymbol{\theta})$ and correspondingly f denotes the joint density of observed data \mathbf{x} and latent variables \mathbf{z} (In other words, $\tilde{\mathcal{D}}(q, f)$ refers to the negative of ELBO). Given the relation demonstrated in (20), our objective is to minimize $\mathcal{D}(q)$ by choosing q . Gradient Boosting is one frequently used method for this objective and specifically for this case the rough idea for implementing gradient boosting is by considering perturbation from q to $(1 - \epsilon)q + \epsilon h$. It would be useful to derive the *functional derivative* of $\mathcal{D}(q)$, denoted by $\nabla \mathcal{D}(q)$, a direct derivation with shorthand notation is as following ⁷

$$\begin{aligned}\nabla \mathcal{D}(q) &= \nabla_q \left(q \log \frac{q}{f} \right) \\ &= \log \frac{q}{f} + q \frac{f}{q} \frac{1}{f} = \log \frac{q}{f} + 1\end{aligned}$$

Recall the gradient boosting dynamics we introduced previously,

$$q_t = (1 - \alpha_t) q_{t-1} + \alpha_t h_t$$

⁷ Here we emphasize that q as the function value should be regarded as a univariate variable

and with $\epsilon \rightarrow 0$, Taylor expansion with respect to ϵ yields

$$\begin{aligned}\mathcal{D}((1-\epsilon)q + \epsilon h) &= \mathcal{D}(q + (h-q)\epsilon) \\ &= \mathcal{D}(q) + \epsilon \langle h-q, g \rangle + o(\epsilon^2)\end{aligned}$$

Consequently for a fixed q we just need to choose h to minimize $\langle h, g \rangle = \langle h, \nabla \mathcal{D}(q) \rangle$ and this also suggests that intuitively we need to approximately choose h to match the “negative” direction $\nabla \mathcal{D}(q)$, i.e. $-\nabla \mathcal{D}(q)$ if possible. And with the derived $\nabla \mathcal{D}(q)$ we focus on the following optimization problem instead

$$\mathcal{D}(q_t) = \mathcal{D}(q_{t-1}) + \alpha_t \langle h_t, \log(q_{t-1}/f) \rangle - \alpha_t \langle q_{t-1}, \log(q_{t-1}/f) \rangle + o(\alpha_t^2) \quad (45)$$

5.2.2 Application of variational bayes method

Recently there are some emerging literature discussing the application of variational expectation maximization method in estimating time-varying parameter models (TVP) (Kowal et al., 2019; Koop and Korobilis, 2020, henceforth [KK20](#)) based on the recent progress (Wang and Blei, 2019; Alquier and Ridgway, 2020) in theoretically justifying the asymptotic properties of variational bayes method. Following discussion is mainly about the basic data generating process (DGP) under TVP framework and how variational bayes method is applied for estimating the parameters specifying TVP modelling. Data is assumed to be generated as following

$$y_t = \beta_{1t}x_{1t} + \beta_{2t}x_{2t} + \dots + \beta_{pt}x_{pt} + \sigma_t \varepsilon_t, \quad \varepsilon_t \sim \mathcal{N}(0, 1) \quad (46)$$

$$x_{j,t} \sim \mathcal{N}(0, 1), \quad j = 1, \dots, p \quad (47)$$

$$\beta_{j,t} = s_{j,t} \times \theta_{j,t}, \quad s_{j,t} \text{ is 0-1 indicator variable} \quad (48)$$

$$\theta_{j,t} = \underline{\theta}_j + \underline{\rho}(\theta_{j,t-1} - \underline{\theta}_j) + \underline{\delta}\eta_{j,t}, \quad \eta_{j,t} \sim \mathcal{N}(0, 1) \quad (49)$$

$$\log(\sigma_t^2) = \underline{\sigma}^2 + \underline{\phi}(\log(\sigma_{t-1}^2) - \underline{\sigma}^2) + \underline{\xi}\zeta_t, \quad \zeta_t \sim \mathcal{N}(0, 1) \quad (50)$$

$$\theta_{j,0} = \underline{\theta}_j, \quad \log(\sigma_0^2) = \underline{\sigma}^2 \quad (51)$$

Generally the model specified under the time-varying parameter setting is as following,

$$y_t = \mathbf{X}_t \boldsymbol{\beta}_t + \epsilon_t \quad (52)$$

$$\boldsymbol{\beta}_t = \boldsymbol{\beta}_{t-1} + \boldsymbol{\eta}_t \quad (53)$$

where $\boldsymbol{\beta}_t = (\beta_{1t}, \dots, \beta_{pt})^\top$; $\epsilon_t \sim \mathcal{N}(0, \sigma_t^2)$ with σ_t^2 as time-varying parameter; $\boldsymbol{\eta}_t \sim \mathcal{N}(0, \mathbf{W}_t)$ with \mathbf{W}_t as diagonal matrix $\mathbf{W}_t = \text{diag}(w_{1t}, \dots, w_{pt})$. For later description simplicity, we introduce the notation $\mathbf{w}_t = [w_{1t}, \dots, w_{pt}]^\top$ as the $p \times 1$ vector collecting diagonal elements of \mathbf{W}_t . To see how variational Bayes method is applied in this time-varying framework. It is the generally the interests

of researchers to recover $p(\boldsymbol{\beta}_t, \mathbf{w}_t | \mathbf{y}_{1:t})$ from Bayesian perspective. We further assume the factorized structure for the variational density

$$q(\boldsymbol{\beta}_t, \mathbf{w}_t | \mathbf{y}_{1:t}) = q(\boldsymbol{\beta}_t | \mathbf{y}_{1:t}) \prod_{j=1}^p q(w_{j,t} | \mathbf{y}_{1:t}) \quad (54)$$

Similar logic applies to get the optimized functional form of ELBO akin to the the discussion of generic variation bayes method, which yields that ELBO is maximized by iterating through the following recursions

$$\begin{aligned} q(\boldsymbol{\beta}_t | \mathbf{y}_{1:t}) &\propto \exp\left(\int \log p(\boldsymbol{\beta}_t, \mathbf{w}_t | \mathbf{y}_{1:t}) \prod_j^p q(w_{j,t} | \mathbf{y}_{1:t}) d\mathbf{w}_t\right) \\ &\propto \exp\left(\int \log p(y_t, \boldsymbol{\beta}_t, \mathbf{w}_t | \mathbf{y}_{1:t-1}) \prod_j^p q(w_{j,t} | \mathbf{y}_{1:t}) d\mathbf{w}_t\right) \end{aligned} \quad (55)$$

$$\begin{aligned} q(w_{j,t} | \mathbf{y}_{1:t}) &\propto \exp\left(\int \log p(\boldsymbol{\beta}_t, \mathbf{w}_t | \mathbf{y}_{1:t}) q(\boldsymbol{\beta}_t | \mathbf{y}_{1:t}) d\boldsymbol{\beta}_t\right) \\ &\propto \exp\left(\int \log p(y_t, \boldsymbol{\beta}_t, \mathbf{w}_t | \mathbf{y}_{1:t-1}) q(\boldsymbol{\beta}_t | \mathbf{y}_{1:t}) d\boldsymbol{\beta}_t\right), \quad j = 1, \dots, p \end{aligned} \quad (56)$$

Then the following recursive relationship can be set up

$$\begin{aligned} q(\boldsymbol{\beta}_t | \mathbf{y}_{1:t}) &\propto \exp\left[\mathbb{E}_{q(\mathbf{w}_t | \mathbf{y}_{1:t})}(\log p(y_t, \boldsymbol{\beta}_t, \mathbf{w}_t | \mathbf{y}_{1:t-1}))\right] \\ &= \exp\left\{\mathbb{E}_{q(\mathbf{w}_t | \mathbf{y}_{1:t})} \log [p(y_t | \boldsymbol{\beta}_t, \mathbf{w}_t) p(\boldsymbol{\beta}_t | \mathbf{y}_{1:t-1}) p(\mathbf{w}_t | \mathbf{y}_{1:t-1})]\right\} \\ &= p(y_t | \boldsymbol{\beta}_t, \mathbf{w}_t) \exp\left\{\mathbb{E}_{q(\mathbf{w}_t | \mathbf{y}_{1:t})} [\log p(\boldsymbol{\beta}_t | \boldsymbol{\beta}_{t-1}) + \log q(\boldsymbol{\beta}_t | \mathbf{y}_{1:t-1})]\right\} \\ &\quad \times \underbrace{\exp\left\{\mathbb{E}_{q(\mathbf{w}_t | \mathbf{y}_{1:t})} [\log p(\mathbf{w}_t | \mathbf{y}_{1:t-1})]\right\} \times [\log q(\boldsymbol{\beta}_t | \mathbf{y}_{1:t-1})]^{-1}}_{\propto \text{const}} \end{aligned} \quad (57)$$

As suggested by George and McCulloch (1993), variable selection under this high-dimensional and dynamic setting is implemented based on the following spike-and-slab prior specification on the time-varying coefficients,

$$\beta_{j,t} | \gamma_{j,t}, \tau_{j,t}^2 \sim (1 - \gamma_{j,t}) \mathcal{N}(0, \underline{c}\tau_{j,t}^2) + \gamma_{j,t} \mathcal{N}(0, \tau_{j,t}^2) \quad (58)$$

$$\gamma_{j,t} | \pi_{0,t} \sim \text{Bernoulli}(\pi_{0,t}) \quad (59)$$

$$\frac{1}{\tau_{j,t}^2} \sim \text{Gamma}(g_0, h_0) \quad (60)$$

$$\pi_{0,t} \sim \text{Beta}(1, 1). \quad (61)$$

That is , If $\gamma_{j,t} = 1$, the prior for $\beta_{j,t}$ has a normal prior with zero mean and variance $\tau_{j,t}^2$, while if $\gamma_{j,t} = 0$ the prior variance becomes $\underline{c}\tau_{j,t}^2$ which also implies that whenever variable selection is

required \underline{c} has to be set as $\underline{c} \rightarrow 0$. But along with (58), we also have dynamic transition such that

$$\beta_{j,t} \mid \beta_{j,t-1}, w_{j,t} \sim \mathcal{N}(\beta_{j,t-1}, w_{j,t}) \quad (62)$$

Then it is possible for us to combine the prior information imposed on $\beta_{j,t}$ jointly with (58) and (62) such that

$$\beta_t = \tilde{\mathbf{F}}_t \beta_{t-1} + \tilde{\boldsymbol{\eta}}_t \quad (63)$$

where

$$\begin{aligned} \tilde{\boldsymbol{\eta}}_t &\sim \mathcal{N}(\mathbf{0}, \tilde{\mathbf{W}}_t) \\ \tilde{\mathbf{W}}_t &= [\mathbb{E}(\mathbf{W}_t)^{-1} + \mathbb{E}(\mathbf{V}_t)^{-1}]^{-1} \\ \tilde{\mathbf{F}}_t &= \tilde{\mathbf{W}}_t \times \mathbb{E}(\mathbf{W}_t)^{-1} \\ \mathbf{W}_t &= \text{diag}(w_{1,t}, \dots, w_{p,t}) \\ \mathbf{V}_t &= \text{diag}(v_{1,t}, \dots, v_{p,t}) \end{aligned}$$

A quick demonstration for the replicated Monte Carlo experiment of KK20 is as following

[Place Figure 2 about here]

[Place Figure 3 about here]

My implementation instead is based on the optimized hybrid R and C++ codes relying on [Microsoft R Open distribution](#). This implementation automatically implements parallel matrix computation based on the C++ Armadillo linear algebra template and consequently would improve computational efficiency significantly on multi-cores system. Armadillo as the C++ library for linear algebra and & scientific computing is initially developed and actively maintained by Conrad Sanderson from Griffith University. More details about hybrid coding in R and C++ are covered comprehensively in the classic textbook Eddelbuettel (2013). The associated R package `vbdvsarmadillo` and more details about implementation are both available at

<https://www.yaohanchen.com/post/vbdvsr/cpp/>

5.3 Variational Bayes with Intractable Likelihood

As in Tran et al. (2017), let us consider the a generic scenario where there are observed random variables \mathbf{y} , latent variables \mathbf{z} and the associated parameters $\boldsymbol{\theta} \in \Theta$. Within the Bayesian analytical framework, in general we are interested in

$$\pi(\boldsymbol{\theta}) = p(\boldsymbol{\theta} \mid \mathbf{y}) \propto p(\mathbf{y} \mid \boldsymbol{\theta})p(\boldsymbol{\theta}) \quad (64)$$

which is usually referred to the posterior distribution on which we make posterior inference. However, $p(\mathbf{y} | \boldsymbol{\theta})$ is required for constructing sampler like MCMC, which under some circumstance is not exactly available. Hence variational bayes method instead provides an alternative for recovering $\pi(\boldsymbol{\theta})$ when the the likelihood function $p(\mathbf{y} | \boldsymbol{\theta})$ is not available. The key idea is discussed as following. In fact, we can firstly apply filtering techniques like particle filter to estimate the likelihood function $p(\mathbf{y} | \boldsymbol{\theta})$, denoted by $\widehat{p}_N(\mathbf{y} | \boldsymbol{\theta})$ with N specifying the number of particles used for estimating, and we use \mathbf{z} to denote the difference between the estimated likelihood function and the true likelihood function in log scale such that

$$\mathbf{z} = \log \widehat{p}_N(\mathbf{y} | \boldsymbol{\theta}) - \log p(\mathbf{y} | \boldsymbol{\theta})$$

which is the form that we can alternatively write as $\widehat{p}_N(\mathbf{y} | \boldsymbol{\theta}) = e^{\mathbf{z}} p(\mathbf{y} | \boldsymbol{\theta})$. Given the unbiasedness of this estimated likelihood function, we thus have the following fact

$$\int e^{\mathbf{z}} g_N(\mathbf{z} | \boldsymbol{\theta}) d\mathbf{z} = 1$$

where $g_N(\mathbf{z} | \boldsymbol{\theta})$ refers to the density function of random variable z conditional on parameter $\boldsymbol{\theta}$. With this newly introduced random variable z , we can define the following augmented likelihood function $\pi_N(\boldsymbol{\theta}, \mathbf{z})$ on $\Theta \times \mathbf{R}$,⁸

$$\pi_N(\boldsymbol{\theta}, \mathbf{z}) = \frac{p(\boldsymbol{\theta}) p(\mathbf{y} | \boldsymbol{\theta}) e^{\mathbf{z}} g_N(\mathbf{z} | \boldsymbol{\theta})}{p(\mathbf{y})} = \pi(\boldsymbol{\theta}) e^{\mathbf{z}} g_N(\mathbf{z} | \boldsymbol{\theta}). \quad (65)$$

From which we can derive that

$$\begin{aligned} \log p(\mathbf{y}) &= \log \left[\frac{p(\boldsymbol{\theta}) P(\mathbf{y} | \boldsymbol{\theta}) e^{\mathbf{z}} g_N(\mathbf{z} | \boldsymbol{\theta})}{\pi_N(\boldsymbol{\theta}, \mathbf{z})} \right] \\ &= \log \left[\frac{p(\boldsymbol{\theta}) \widehat{p}_N(\mathbf{y} | \boldsymbol{\theta})}{q_{\lambda, N}(\boldsymbol{\theta}, \mathbf{z})} \right] + \log \left[\frac{q_{\lambda, N}(\boldsymbol{\theta}, \mathbf{z}) g_N(\boldsymbol{\theta}, \mathbf{z})}{\pi_N(\boldsymbol{\theta}, \mathbf{z})} \right] \end{aligned} \quad (66)$$

where $q_{\lambda, N}(\boldsymbol{\theta}, \mathbf{z})$ is the correspondingly introduced density function over the space of $\boldsymbol{\theta}$ and \mathbf{z} . Hence taking integration on both sides over $\boldsymbol{\theta}, \mathbf{z}$ with respect to $q_{\lambda, N}(\boldsymbol{\theta}, \mathbf{z})$ on both sides of the above

⁸ Augmented likelihood function defined here is guaranteed to be well-defined density function given that $\int e^{\mathbf{z}} g_N(\mathbf{z} | \boldsymbol{\theta}) d\mathbf{z} = 1$ and $\pi(\boldsymbol{\theta})$ is well-defined posterior. In other words, integration of $\pi_N(\boldsymbol{\theta}, \mathbf{z})$ over $\boldsymbol{\theta}$ and \mathbf{z} is exactly equal 1.

equation yields

$$\log p(\mathbf{y}) = \underbrace{\int q_{\lambda,N}(\boldsymbol{\theta}, \mathbf{z}) \log \left[\frac{p(\boldsymbol{\theta}) \widehat{p}_N(\mathbf{z} | \boldsymbol{\theta})}{q_{\lambda,N}(\boldsymbol{\theta}, \mathbf{z})} \right] d\boldsymbol{\theta} d\mathbf{z}}_{\text{LB}(\lambda)} + \underbrace{\int q_{\lambda,N}(\boldsymbol{\theta}, \mathbf{z}) \log \left[\frac{q_{\lambda,N}(\boldsymbol{\theta}, \mathbf{z}) g_N(\boldsymbol{\theta}, \mathbf{z})}{\pi_N(\boldsymbol{\theta}, \mathbf{z})} \right] d\boldsymbol{\theta} d\mathbf{z}}_{\text{KL}(\lambda)} \quad (67)$$

Remark 5.5

- $q_{\lambda,N}(\boldsymbol{\theta}, \mathbf{z})$ as the constructed density function can especially take the following form

$$q_{\lambda,N}(\boldsymbol{\theta}, \mathbf{z}) = q_{\lambda}(\boldsymbol{\theta}) g_N(\boldsymbol{\theta}, \mathbf{z})$$

and this structure separates the dependence of λ on N and hence could bring some convenience, this will be discussed later.

- $\text{KL}(\lambda)$ is the associated Kullback-Leibler divergence by comparing $q_{\lambda,N}(\boldsymbol{\theta}, \mathbf{z})$ with the true augmented likelihood function $\pi_N(\boldsymbol{\theta}, \mathbf{z})$. This is a function of λ and always no less than 0.
- $\text{LB}(\lambda)$ as the function of λ is justified as the lower bound of the marginal log likelihood $\log p(\mathbf{y})$ given the fact that $\text{KL}(\lambda) \geq 0$. $\text{LB}(\lambda)$ is later applied for checking the convergence of gradient.

5.4 Dynamic Shrinking Process

Before we move on to the application of variational bayes method, it necessary to briefly discuss a related branch of literature recently appearing in statistics discussing shrinkage dynamically. One of the representative work is done by Kowal et al. (2019), and the following discussion is mainly based on the framework of this paper. All the stuff to be discussed is closely related with the concepts about global-local prior defined as following

$$\omega_t | \tau, \lambda_t \overset{\text{indep}}{\sim} \mathcal{N}(0, \tau^2 \lambda_t^2), \quad (68)$$

And $h_t = \log(\tau^2 \lambda_t^2)$ is a process following the general dependent data model

$$h_t = \mu + \psi_t + \eta_t, \quad \eta_t \overset{\text{iid}}{\sim} Z(\alpha, \beta, 0, 1) \quad (69)$$

where

$$\begin{aligned} \mu &= \log(\tau^2) \\ \psi_t + \eta_t &= \log(\lambda_t^2) \end{aligned}$$

and $Z(\alpha, \beta, \mu_z, \sigma_z)$ denotes the Z -distribution with density function specified as following

$$[z] = [\sigma_z B(\alpha, \beta)]^{-1} \{ \exp[(z - \mu_z) / \sigma_z] \}^\alpha \{ 1 + \exp[(z - \mu_z) / \sigma_z] \}^{-(\alpha + \beta)}, z \in \mathbb{R} \quad (70)$$

and $B(\cdot, \cdot)$ refers to the Beta function. Then the dynamic shrinkage process is modelled as following

$$h_{t+1} = \mu + \phi(h_t - \mu) + \eta_t, \quad \eta_t \stackrel{iid}{\sim} Z(\alpha, \beta, 0, 1) \quad (71)$$

Remark 5.6 ψ_t plays the role governing the adaptive shrinkage process with the following two emphasized special case

- Replacing ψ_t in (70) with $\phi(h_t - \mu)$ in (71) yields equivalence.
- Replacing ψ_t in (70) with $\mathbf{z}_t^\top \boldsymbol{\alpha}$ for a vector of predictors yields the equivalence for linear regression framework.

Remark 5.7 Z -distribution arises from Beta distribution by the following discussion

- $B(\alpha, \beta) = \int_0^1 u^{\alpha-1} (1-u)^{\beta-1} du = \int_0^1 (1-u)^{\alpha-1} u^{\beta-1} du = B(\beta, \alpha)$.
- Define $u = \frac{1}{1 + \exp[(z - \mu_z) / \sigma_z]}$, we can readily check the $[z]$ of the previous functional form is indeed a probability density function,

$$\begin{aligned} & \int [\sigma_z B(\alpha, \beta)]^{-1} \{ \exp[(z - \mu_z) / \sigma_z] \}^\alpha \{ 1 + \exp[(z - \mu_z) / \sigma_z] \}^{-(\alpha + \beta)} dz \\ &= \frac{1}{\sigma_z B(\alpha, \beta)} \int_0^1 (1-u)^\alpha u^\beta \left| \frac{-\sigma_z}{(1-u)u} \right| du = \frac{\sigma_z}{B(\alpha, \beta)} \int_0^1 (1-u)^{\alpha-1} u^{\beta-1} du = 1 \end{aligned}$$

- The first moment of Z -distribution is μ_z ,

$$\begin{aligned} & \int [\sigma_z B(\alpha, \beta)]^{-1} \{ \exp[(z - \mu_z) / \sigma_z] \}^\alpha \{ 1 + \exp[(z - \mu_z) / \sigma_z] \}^{-(\alpha + \beta)} z dz \\ &= \frac{1}{\sigma_z B(\alpha, \beta)} \int_0^1 (1-u)^\alpha u^\beta \left[\log\left(\frac{1}{u} - 1\right) \sigma_z + \mu_z \right] \left| \frac{-\sigma_z}{(1-u)u} \right| du = \mu_z \end{aligned}$$

where we use the fact that

$$\begin{aligned}\frac{1}{B(\alpha, \beta)} \int_0^1 (1-u)^{\alpha-1} u^{\beta-1} \log(1-u) du &= \frac{1}{B(\alpha, \beta)} \int_0^1 u^{\alpha-1} (1-u)^{\beta-1} \log u du \\ &= \psi(\alpha) - \psi(\alpha + \beta) \\ \frac{1}{B(\alpha, \beta)} \int_0^1 (1-u)^{\alpha-1} u^{\beta-1} \log u du &= \psi(\beta) - \psi(\beta + \alpha)\end{aligned}$$

which implies that the first moment of $Z(\alpha, \beta, \mu_z, \sigma_z)$ should be $\psi(\alpha) - \psi(\beta) + \mu_z$ and hence it is symmetric with the origin concentrating on μ_z if and only if $\alpha = \beta$. $\psi(\cdot)$ refers to digamma function, which is essentially the logarithm derivative of gamma function.

Remark 5.8 Some useful discussion about special distributions

- If specifically a random variable is following Z distribution such that $\eta \sim Z(\alpha, \beta, \mu_z, 1)$, $\mu_z \in \mathbb{R}$. Then $\kappa = 1/(1 + \exp(\eta)) \sim TPB(\beta, \alpha, \exp(\mu_z))$, where $\kappa \sim TPB(\beta, \alpha, \gamma)$ denotes the three parameter Beta distribution (Armagan et al., 2011) with density specified as following

$$[\kappa] = [B(\beta, \alpha)]^{-1} \gamma^\beta \kappa^{\beta-1} (1-\kappa)^{\alpha-1} [1 + (\gamma-1)\kappa]^{-(\alpha+\beta)}, \kappa \in (0, 1), \gamma > 0 \quad (72)$$

Proof. Recall that $\eta \sim Z(\alpha, \beta, \mu_z, 1)$ has the following distribution with $\sigma_z = 1$

$$[z] = [\sigma_z B(\alpha, \beta)]^{-1} \{\exp[(z - \mu_z)/\sigma_z]\}^\alpha \{1 + \exp[(z - \mu_z)/\sigma_z]\}^{-(\alpha+\beta)}$$

which implies that if we define $\lambda^2 = \exp(\eta)$, then

$$\begin{aligned}[\lambda^2] &\propto (\lambda^2)^{-1} \{\exp[\log(\lambda^2) - \mu_z]\}^\alpha \{1 + \exp[\log(\lambda^2) - \mu_z]\}^{-(\alpha+\beta)} \\ &\propto (\lambda^2)^{\alpha-1} [1 + \lambda^2/\exp(\mu_z)]^{-(\alpha+\beta)}\end{aligned}$$

and for $\kappa = 1/(1 + \lambda^2)$,

$$\begin{aligned}[\kappa] &\propto \kappa^{-2} [\kappa^{-1} - 1]^{\alpha-1} [1 + (\kappa^{-1} - 1)/\exp(\mu_z)]^{-(\alpha+\beta)} \\ &\propto \kappa^{-2-(\alpha-1)} (1-\kappa)^{\alpha-1} \{\kappa^{-1} [\kappa \exp(\mu_z) + (1-\kappa)]\}^{-(\alpha+\beta)} \\ &\propto (1-\kappa)^{\alpha-1} \kappa^{\beta-1} [\kappa \exp(\mu_z) + (1-\kappa)]^{-(\alpha+\beta)}\end{aligned}$$

□

- In a general framework, we consider λ^2 distributed with density

$$[\lambda^2] \propto (\lambda^2)^{\alpha-1} (1 + \lambda^2)^{-(\alpha+\beta)}, \lambda > 0 \quad (73)$$

which is referred to inverted-Beta distribution, $IB(\beta, \alpha)$. Similar to the previous discussion, we can show that $\kappa = 1/(1 + \lambda^2) \sim Beta(\beta, \alpha)$.

Proof.

$$\begin{aligned} [\kappa] &\propto \kappa^{-2} \left(\frac{1}{\kappa} - 1\right)^{\alpha-1} \left(\frac{1}{\kappa}\right)^{-(\alpha+\beta)} \\ &\propto (1 - \kappa)^{\alpha-1} \kappa^{-2+1-\alpha+\alpha+\beta} = (1 - \kappa)^{\alpha-1} \kappa^{\beta-1}, \end{aligned}$$

□

- Previous discussion hence can be summarized as following claim

$$\lambda^2 \sim IB(\beta, \alpha) \Leftrightarrow \kappa = 1/(1 + \lambda^2) \sim Beta(\beta, \alpha) \Leftrightarrow \eta = \log(\lambda^2) = \log(\kappa^{-1} - 1) \sim Z(\alpha, \beta, 0, 1).$$

Let us temporarily focus on the following process and discuss some associated properties

$$\kappa_{t+1} = \frac{1}{1 + \tau^2 \lambda_{t+1}^2} \quad (74)$$

The following theorem plays the pivotal rule

Theorem 5.3 *For the dynamic process specified as in (71), the conditional prior distribution of κ_{t+1} is*

$$[\kappa_{t+1} \mid \{\kappa_s\}_{s \leq t}, \phi, \tau] \sim TPB\left(\beta, \alpha, \tau^{2(1-\phi)} \left[\frac{1 - \kappa_t}{\kappa_t}\right]^\phi\right) \quad (75)$$

Proof. Recall the dynamic process specified as in (71), we have

$$[h_{t+1} \mid h_t, \phi, \mu] \sim Z(\alpha, \beta, \mu + \phi(h_t - \mu), 1)$$

By using the results discussed in the aforementioned remarks, we have the conditional distribution of κ_{t+1} is

$$[\kappa_{t+1} \mid h_t, \phi, \mu] \sim TPB(\beta, \alpha, \exp(\mu + \phi(h_t - \mu))) \quad (76)$$

□

Theorem 5.4 *For the dynamic process specified as in (71) and $\alpha = \beta = \frac{1}{2}$, the conditional prior distribution satisfies*

$$\mathbb{P}(\kappa_{t+1} < \varepsilon \mid \{\kappa_s\}_{s \leq t}, \phi) \rightarrow 1 \quad (77)$$

as $\kappa_t \rightarrow 0$ or any fixed $\varepsilon \in (0, 1)$ and $\phi \neq 0$

It is worthwhile discussing Pólya-Gamma distribution for introducing techniques mentioned later that have widely discussed in literature such as Barndorff-Nielsen et al. (1982) and Polson et al. (2013). The Pólya-Gamma distributions, denoted by $PG(b, c)$ is a subset of the class of infinite

convolutions of gamma distributions. As a special case, $PG(1,0)$ is referred to the carefully chosen element of the the class of infinite convolutions of exponentials. In general, we have the following factorization, all the listed results are the consequence of Weierstrass factorization theorem

$$\begin{aligned}\sin(\pi z) &= \pi z \prod_{k=1}^{\infty} \left(1 - \left(\frac{1}{k}\right)^2\right) \\ \cos(\pi z) &= \prod_{k=0}^{\infty} \left(1 - \left(\frac{z}{k + \frac{1}{2}}\right)^2\right) \\ \sinh(z) &= z \prod_{k=1}^{\infty} \left(1 + \frac{z^2}{(\pi k)^2}\right) \\ \cosh(z) &= \prod_{k=1}^{\infty} \left(1 + \frac{z^2}{(\pi k - \pi/2)^2}\right)\end{aligned}$$

Based on these facts implied from Weierstrass factorization and that

- If a random variable ω is distributed following $PG(1,0)$, we are able to obtain the Laplace transform (with characterization function evaluated at it , where i denotes imaginary unit) as following

$$\mathbb{E}[\exp(-\omega t)] = \cosh^{-1}(\sqrt{t/2})$$

Moreover, if ω is distributed following $PG(b,0)$, the corresponding Laplace transform is as following

$$\mathbb{E}[\exp(-\omega t)] = \prod_{k=1}^{\infty} \left(1 + \frac{t}{2\pi^2(k-1/2)^2}\right)^{-b} = \frac{1}{\cosh^b(\sqrt{t/2})}. \quad (78)$$

- Gamma distribution denoted as $\text{Ga}(\alpha, \beta)$, the characterization function generally takes the following form

$$\left(1 - \frac{i\tau}{\beta}\right)^{-\alpha}$$

and accordingly with τ evaluated at it (complex number with only imaginary part), we have Laplace transform of Gamma distribution as following

$$\left(1 + \frac{t}{\beta}\right)^{-\alpha}$$

Replacing $\alpha = b$ and $\beta = 1$ directly implies that Laplace transform of $\text{Ga}(b,1)$ is

$$(1+t)^{-b}$$

- Suppose that random variable ω is constructed as following

$$\omega \stackrel{D}{=} \frac{1}{2\pi^2} \sum_{k=1}^{\infty} \frac{g_k}{(k-1/2)^2}$$

and $g_k \sim \text{Ga}(b, 1)$ (g_k are mutually independent), then given that the Laplace transform of $\text{Ga}(b, 1)$ is $(1+t)^{-1}$, we have the Laplace transform of this constructed random variable ω takes exactly the form of (78). This is why ω refers to Pólya-Gamma distribution.

As for the extension from $PG(b, 0)$ to $PG(b, c)$, just introduce the following density function for $PG(b, c)$.

$$p(\omega | b, c) = \frac{\exp\left(-\frac{c^2}{2}\omega\right) p(\omega | b, 0)}{\mathbb{E}_{\omega} \left\{ \exp\left(-\frac{c^2}{2}\omega\right) \right\}} \quad (79)$$

where \mathbb{E}_{ω} is taken with respect to $p(\omega | b, 0)$. Based on the detailed discussion in Polson et al. (2013) (equation (6) specifically), it is possible to demonstrate that $\omega \sim PG(b, c)$ is equivalent as following

$$\omega \stackrel{D}{=} \sum_{k=1}^{\infty} \frac{\text{Ga}(b, 1)}{d_k} = \frac{1}{2\pi^2} \sum_{k=1}^{\infty} \frac{\text{Ga}(b, 1)}{\left(k - \frac{1}{2}\right)^2 + c^2 / (4\pi^2)} \quad (80)$$

5.5 Quick note EB coverage interval

Useful facts summarized from (Morris, 1983; Carlin and Louis, 2000), If we specify likelihood as normal such that $f(y | \theta) = \frac{1}{\sigma\sqrt{2\pi}} \exp\left(-\frac{(y-\theta)^2}{2\sigma^2}\right)$, and $\theta \sim \mathcal{N}(\mu, \tau^2)$, $\theta \in \mathbf{R}$, $\mu \in \mathbf{R}$. Furthermore, we assume that σ is a constant positive known parameter and μ, τ are known hyperparameters as well. The posterior of θ for this specification is

$$p(\theta | y) = \mathcal{N}\left(\theta \mid \frac{\sigma^2\mu + \tau^2y}{\sigma^2 + \tau^2}, \frac{\sigma^2\tau^2}{\sigma^2 + \tau^2}\right)$$

For a specific setting ($\mu = 0$), EB estimator for θ_i is given as $w_{EB}Y_i$ with $w_{EB} = \frac{\tau^2}{\sigma^2 + \tau^2}$. we want to check whether θ_i is contained in a posterior stochastic interval for a specific critical value χ (in some sense, χ could be interpreted as radius of neighbourhood with its center on $w_{EB}Y_i$). Or in other words, non-coverage posterior confidence interval refers to those stochastic intervals, CI, $w_{EB}Y_i \pm \chi \cdot w_{EB}\sigma$ in which a fixed θ_i is not contained. All these intervals correspond to the case

$$|\theta_i - w_{EB}Y_i| \geq \chi \cdot w_{EB}\sigma \Leftrightarrow \left| \frac{\theta_i - w_{EB}Y_i}{w_{EB}\sigma} \right| \geq \chi$$

But note that

$$\mathbb{E}[w_{EB}Y_i - \theta_i] = (w_{EB} - 1)\theta_i \quad \text{Var}[w_{EB}Y_i - \theta_i] = w_{EB}^2\sigma^2$$

hence the corresponding defined statistics $\frac{w_{EB}Y_i - \theta_i}{w_{EB}\sigma}$ is normally distributed with mean equal to b_i

and variance equal to 1, where b_i is of the following form

$$b_i = \frac{(w_{EB} - 1)\theta_i}{w_{EB}\sigma}$$

which justifies the definition such that

$$\frac{\theta_i - w_{EB}Y_i}{w_{EB}\sigma} - (-b_i) = Z \sim \mathcal{N}(0, 1) \text{ and } \left| \frac{\theta_i - w_{EB}Y_i}{w_{EB}\sigma} \right| = |Z - b_i|$$

and the claim that probability of non-coverage is

$$r(b_i, \chi) = \mathbb{P}\left(\left| \frac{\theta_i - w_{EB}Y_i}{w_{EB}\sigma} \right| \geq \chi\right) = \mathbb{P}(|Z - b_i| \geq \chi) = \Phi(-\chi - b) + \Phi(-\chi + b) \quad (81)$$

where $\Phi(\cdot)$ denotes the cdf of standard normal distribution. $\Phi(\cdot)$ is differentiable and hence $r(b_i, \chi)$ is differentiable w.r.t. b_i , which is important for constructing the optimal functional form of $r(b, \chi)$. Note for the all the previous discussion, we keep our discussion based on fixing θ_i . Under Bayesian framework or the scenario where random effect is allowed, b_i by definition as the function of θ_i inherits the randomness from θ_i . In general θ_i can follow any form of feasible distribution but we obtain our EB estimator by pretending imposing normality. As the result, once we consider integrating out the uncertainty from b_i , which is actually inherited from θ_i , we just impose the structure that $b_i \sim F$, where F is used for generally describing the distribution of b_i .

Remark 5.9

- Given that Z as the standard normal distribution is symmetric on \mathbf{R} , $r(b, \chi)$ is symmetric in b .
- For all $t \geq 0$, we can define $r_0(t, \chi) = r(\sqrt{t}, \chi)$ and hence the following two optimization problems are equivalent

$$\begin{aligned} \sup_F \mathbb{E}_F[r(t, \chi)] & \Leftrightarrow \sup_F \mathbb{E}_F[r_0(t, \chi)] \\ \text{s. t. } \mathbb{E}_F[t^2] = m_2 & \qquad \text{s. t. } \mathbb{E}_F[t] = m_2 \end{aligned} \quad (82)$$

and we denote the corresponding optimal value as $\rho(m_2, \chi)$.

- The following statement (or lemma) discussed in Carolan (2002) gives a relatively formal definition about *least concave majorant* for the univariate case.

Claim: Suppose g is a function defined on a set containing $[a, b]$, where a can equal $-\infty$ or b can equal ∞ . Then g^* defined as the least concave majorant of g for $t \in [a, b]$ is

$$g^*(t) = \sup_{a \leq x_1 \leq t} \sup_{t \leq x_2 \leq b} \left\{ \frac{(x_2 - t)g(x_1) + (t - x_1)g(x_2)}{x_2 - x_1} \right\} \quad (\dagger)$$

where $0/0$ is defined as $g(t)$ when $x_1 = x_2 = t$. Alternatively, we can concisely define g^* as

$$g^*(t) = \inf \{ \tilde{g}(t) : \tilde{g}_t \geq g(t), \tilde{g} \text{ concave} \}.$$

- Let $\bar{r}(t, \chi)$ be the *least concave majorant* of $r_0(t, \chi)$, it can be shown that $\rho(m_2, \chi) = \bar{r}(m_2, \chi)$, where $\rho(m_2, \chi)$ refers to the optimal functional form that we want to obtain when t is evaluated at $t = m_2$.
- Previous bullet point implies that optimal functional form of $\rho(t, \chi)$ with t evaluated at the moment restriction $t = m_2$ is essentially given by the least concave majorant of $r_0(t, \chi)$, $\bar{r}(t, \chi)$. A natural question then will be what is the exactly the functional form of $\bar{r}(t, \chi)$. By the definition of least concave majorant, a straightforward motivating construction will be using linear combination as the approximation. It can be proved that $\bar{r}(t, \chi)$ takes the following functional form akin to (†)

$$\bar{r}(t, \chi) = \sup_{u \geq t} \left\{ (1 - t/u)r_0(0, \chi) + \frac{t}{u}r_0(u, \chi) \right\}, \quad (83)$$

6 Latent Factors

Dimension reduction is not only the major concern in machine learning literature but also other fields like Finance and Statistics. Factor structure as the major dimension reduction tools have been widely discussed in literature, both from theoretical and empirical perspective. Recently there are some meaningful discussions (see Lettu and Pelger, 2020a,b) about extending the conventional methodologies of principal component analysis (PCA) to a relatively more general framework by constructing the objective function accounting for the factor approximation error simultaneously from the time-series and cross-section dimension.

Basically, we observe excess return of N assets over T periods and factor structure is imposed as following

$$X_{t,n} = F_t^\top \Lambda_n + e_{t,n} \quad n = 1, \dots, N \quad t = 1, \dots, T \quad (84)$$

where if there are K factors, then both F_t and Λ_n are a $K \times 1$ vectors. In matrix notation this reads as

$$\underbrace{\mathbf{X}}_{T \times N} = \underbrace{\mathbf{F}}_{T \times K} \underbrace{\mathbf{\Lambda}^\top}_{K \times N} + \underbrace{\mathbf{e}}_{T \times N} \quad (85)$$

In practice, data is usually demeaned firstly and then applied with PCA. Moreover since Factor structure is imposed for summarizing the time-series variation shared commonly by cross-sectional data, the demeaned process is taken over time-series dimension. For description simplicity, introducing

the following notation first,

$$\text{projection matrix demeaning time-series : } \mathbf{M}_1 = \mathbf{I}_T - \frac{1}{T}\iota\iota^\top$$

$$\text{projection matrix : } \mathbf{M}_\Lambda = \mathbf{I}_N - \Lambda (\Lambda^\top \Lambda)^{-1} \Lambda^\top$$

where ι is a $T \times 1$ vector of 1's. With this notation introduced, the demeaned data (across time-series dimension) is readily represented in matrix notation as

$$\tilde{\mathbf{X}} = \mathbf{M}_1 \mathbf{X}$$

$$\tilde{\mathbf{F}} = \mathbf{M}_1 \mathbf{F}$$

The key step in implementing PCA with respect to $\tilde{\mathbf{X}}$ is calculating the sample second moment of $\tilde{\mathbf{X}}$:

$$\frac{1}{T} \tilde{\mathbf{X}}^\top \tilde{\mathbf{X}} = \frac{1}{T} \mathbf{X}^\top \mathbf{M}_1 \mathbf{X} = \frac{1}{T} \mathbf{X}^\top \mathbf{X} - \bar{\mathbf{X}} \bar{\mathbf{X}}^\top$$

Denote $\hat{\Lambda}$, $\hat{\mathbf{F}}$ as the estimated factor loadings and factors respectively, the following equivalence can be proved

Proposition 6.1 *The k -th column of $\hat{\Lambda}$ is proportional to the k -th eigenvector extracted from*

$$\frac{1}{T} \tilde{\mathbf{X}}^\top \tilde{\mathbf{X}} = \frac{1}{T} \mathbf{X}^\top \mathbf{X} - \bar{\mathbf{X}} \bar{\mathbf{X}}^\top$$

where $\bar{\mathbf{X}}$ as the $N \times 1$ vector denotes the time-series sample mean of cross-sectional excess returns, that is

$$\bar{\mathbf{X}} = \frac{1}{T} \mathbf{X}^\top \iota$$

and $\hat{\Lambda}$, $\hat{\mathbf{F}}$ jointly solves the following minimization problem

$$\min_{\Lambda, \mathbf{F}} \frac{1}{NT} \sum_{n=1}^N \sum_{t=1}^T (\tilde{\mathbf{X}}_{t,n} - \tilde{F}_t^\top \Lambda_n)^2 = \min_{\Lambda} \frac{1}{NT} \text{trace} \left[(\tilde{\mathbf{X}} \mathbf{M}_\Lambda)^\top \tilde{\mathbf{X}} \mathbf{M}_\Lambda \right] \quad (86)$$

and

$$\tilde{\mathbf{F}} = \tilde{\mathbf{X}} \Lambda (\Lambda^\top \Lambda)^{-1}$$

Proof. Justification for $\tilde{\mathbf{F}} = \tilde{\mathbf{X}} \Lambda (\Lambda^\top \Lambda)^{-1}$ is simply based on the observation that we are able to commute the summation order of the objective function in a way such that the first summation is taken over n with t fixed and then the second summation is taken over t . First order condition for the summation over n with t fixed is the simply the OLS-style first order condition. Then substituting a given $\tilde{\mathbf{F}}$ of this form back yields the error term formulated as

$$\tilde{\mathbf{X}} - \tilde{\mathbf{F}} \Lambda^\top = \tilde{\mathbf{X}} - \tilde{\mathbf{X}} \Lambda (\Lambda^\top \Lambda)^{-1} \Lambda^\top = \tilde{\mathbf{X}} \mathbf{M}_\Lambda. \quad (87)$$

It is obvious that the objective function of this minimization problem could be represented as

$$\begin{aligned}
& \frac{1}{NT} \text{trace} \left[(\tilde{\mathbf{X}}\mathbf{M}_\Lambda)^\top \tilde{\mathbf{X}}\mathbf{M}_\Lambda \right] \\
&= \frac{1}{NT} \text{trace} \left[\mathbf{M}_\Lambda \tilde{\mathbf{X}}^\top \tilde{\mathbf{X}}\mathbf{M}_\Lambda \right] \\
&= \frac{1}{NT} \text{trace} \left[\tilde{\mathbf{X}}^\top \tilde{\mathbf{X}}\mathbf{M}_\Lambda \right] \\
&= \frac{1}{NT} \text{trace} \left[\tilde{\mathbf{X}}^\top \tilde{\mathbf{X}} \right] - \frac{1}{NT} \text{trace} \left[\tilde{\mathbf{X}}^\top \tilde{\mathbf{X}}\boldsymbol{\Lambda} (\boldsymbol{\Lambda}^\top \boldsymbol{\Lambda})^{-1} \boldsymbol{\Lambda}^\top \right]
\end{aligned}$$

Following discussion is useful for the final proof.

1. $\tilde{\mathbf{X}}^\top \tilde{\mathbf{X}}$ is symmetric real matrix and hence theoretically this is supposed to be represented as

$$\mathbf{Q}\mathbf{U}\mathbf{Q}^\top, \quad \mathbf{U} = \text{diag}(u_1, \dots, u_N), \quad u_1 \geq \dots \geq u_N \geq 0, \quad \mathbf{Q}^\top \mathbf{Q} = \mathbf{Q}\mathbf{Q}^\top = \mathbf{I}_N$$

2. $\boldsymbol{\Lambda} (\boldsymbol{\Lambda}^\top \boldsymbol{\Lambda})^{-1} \boldsymbol{\Lambda}^\top$ is an idempotent matrix, which implies it is supposed to be able to be represented as

$$\mathbf{P}\mathbf{V}\mathbf{P}^\top, \quad \mathbf{V} \text{ is a diagonal matrix with diagonal entries either 1 or 0, } \mathbf{P}^\top \mathbf{P} = \mathbf{P}\mathbf{P}^\top = \mathbf{I}_N$$

These two points jointly imply that

$$\begin{aligned}
\min_{\boldsymbol{\Lambda}} \frac{1}{NT} \text{trace} \left[(\tilde{\mathbf{X}}\mathbf{M}_\Lambda)^\top \tilde{\mathbf{X}}\mathbf{M}_\Lambda \right] &\Leftrightarrow \max_{\boldsymbol{\Lambda}} \frac{1}{NT} \text{trace} \left[\mathbf{Q}\mathbf{U}\mathbf{Q}^\top \mathbf{P}\mathbf{V}\mathbf{P}^\top \right] \\
&= \frac{1}{NT} \text{trace} \left[\mathbf{U}\mathbf{Q}^\top \mathbf{P}\mathbf{V}\mathbf{P}^\top \mathbf{Q} \right]
\end{aligned}$$

Hence the upper bound for $\frac{1}{NT} \text{trace} [\mathbf{U}\mathbf{Q}^\top \mathbf{P}\mathbf{V}\mathbf{P}^\top \mathbf{Q}]$ is the summation of the largest K eigenvalues u_1, \dots, u_K of $\tilde{\mathbf{X}}^\top \tilde{\mathbf{X}}$,⁹ which is necessarily guaranteed by setting the first K columns of \mathbf{Q} proportional to $\boldsymbol{\Lambda}$. \square

However, arbitrage-theory predicts that factors should approximately price the cross-section of expected excess returns well, which also implies that the excess returns should directly enter into the objective function rather than the demeaned counterpart substituted in objective function as in the standard PCA procedure. Hence the following minimization problem should be considered

⁹ A key fact for justifying this claim is $\mathbf{Q}^\top \mathbf{Q} = \mathbf{Q}\mathbf{Q}^\top = \mathbf{I}_N$ and $\mathbf{P}^\top \mathbf{P} = \mathbf{P}\mathbf{P}^\top = \mathbf{I}_K$.

instead.¹⁰

$$\min_{\mathbf{A}, \mathbf{F}} \frac{1}{N} \sum_{n=1}^N \left(\frac{1}{T} \mathbf{X}_n^\top \boldsymbol{\iota} - \mathbf{A}_n^\top \frac{1}{T} \mathbf{F}^\top \boldsymbol{\iota} \right)^2 = \min_{\mathbf{A}} \frac{1}{N} \text{trace} \left[\left(\frac{1}{T} \boldsymbol{\iota}^\top \mathbf{X} \mathbf{M}_{\mathbf{A}} \right) \left(\frac{1}{T} \boldsymbol{\iota}^\top \mathbf{X} \mathbf{M}_{\mathbf{A}} \right)^\top \right] \quad (88)$$

However, this objective function (88) does not identify a set of factors and loadings and the problem admits an infinite number of solutions. In fact, any \mathbf{A} such that $\mathbf{X}^\top \boldsymbol{\iota}$ belongs to the space spanned by the columns of \mathbf{A} will be a solution. A natural step following for extension is by considering the combination of two objective functions in the following way, which has been proposed in Lettu and Pelger (2020a)

$$\begin{aligned} & \min_{\mathbf{A}} \left\{ \frac{1}{NT} \text{trace} \left[\left(\tilde{\mathbf{X}} \mathbf{M}_{\mathbf{A}} \right)^\top \tilde{\mathbf{X}} \mathbf{M}_{\mathbf{A}} \right] + (1 + \gamma) \frac{1}{N} \text{trace} \left[\left(\frac{1}{T} \boldsymbol{\iota}^\top \mathbf{X} \mathbf{M}_{\mathbf{A}} \right) \left(\frac{1}{T} \boldsymbol{\iota}^\top \mathbf{X} \mathbf{M}_{\mathbf{A}} \right)^\top \right] \right\} \quad (89) \\ &= \min_{\mathbf{A}} \left\{ \frac{1}{NT} \text{trace} \left[\left(\mathbf{M}_{\mathbf{A}} \mathbf{X}^\top \mathbf{M}_1 \right) \mathbf{M}_1 \mathbf{X} \mathbf{M}_{\mathbf{A}} \right] + (1 + \gamma) \frac{1}{N} \text{trace} \left[\left(\frac{1}{T} \boldsymbol{\iota}^\top \mathbf{X} \mathbf{M}_{\mathbf{A}} \right) \left(\frac{1}{T} \boldsymbol{\iota}^\top \mathbf{X} \mathbf{M}_{\mathbf{A}} \right)^\top \right] \right\} \\ &= \min_{\mathbf{A}} \left\{ \frac{1}{NT} \text{trace} \left[\mathbf{M}_{\mathbf{A}} \mathbf{X}^\top \left(\mathbf{I}_T - \frac{1}{T} \boldsymbol{\iota} \boldsymbol{\iota}^\top \right) \mathbf{X} \mathbf{M}_{\mathbf{A}} \right] + (1 + \gamma) \frac{1}{N} \text{trace} \left[\left(\frac{1}{T} \mathbf{M}_{\mathbf{A}} \mathbf{X}^\top \boldsymbol{\iota} \right) \left(\frac{1}{T} \boldsymbol{\iota}^\top \mathbf{X} \mathbf{M}_{\mathbf{A}} \right) \right] \right\} \\ &= \min_{\mathbf{A}} \left\{ \frac{1}{NT} \text{trace} \left[\mathbf{M}_{\mathbf{A}} \mathbf{X}^\top \left(\mathbf{I}_T - \frac{1}{T} \boldsymbol{\iota} \boldsymbol{\iota}^\top \right) \mathbf{X} \mathbf{M}_{\mathbf{A}} \right] + (1 + \gamma) \frac{1}{NT} \text{trace} \left[\mathbf{M}_{\mathbf{A}} \mathbf{X}^\top \frac{1}{T} \boldsymbol{\iota} \boldsymbol{\iota}^\top \mathbf{X} \mathbf{M}_{\mathbf{A}} \right] \right\} \\ &= \min_{\mathbf{A}} \frac{1}{NT} \text{trace} \left[\mathbf{M}_{\mathbf{A}} \mathbf{X}^\top \left(\mathbf{I}_T + \frac{\gamma}{T} \boldsymbol{\iota} \boldsymbol{\iota}^\top \right) \mathbf{X} \mathbf{M}_{\mathbf{A}} \right] \quad (90) \end{aligned}$$

The implication from (90) is that the objective function is minimized by choosing \mathbf{A} as the eigenvectors of the largest K eigenvalues of $\mathbf{X}^\top \left(\mathbf{I}_T + \frac{\gamma}{T} \boldsymbol{\iota} \boldsymbol{\iota}^\top \right) \mathbf{X}$.

Remark 6.1 To see how this modified PCA is connected with conventional PCA, we make the following discussion for comparison

	Conventional PCA	Modified PCA
Objective function	$\min_{\mathbf{A}} \frac{1}{NT} \text{trace} \left[\mathbf{M}_{\mathbf{A}} \tilde{\mathbf{X}}^\top \tilde{\mathbf{X}} \mathbf{M}_{\mathbf{A}} \right]$	$\min_{\mathbf{A}} \frac{1}{NT} \text{trace} \left[\mathbf{M}_{\mathbf{A}} \mathbf{X}^\top \left(\mathbf{I}_T + \frac{\gamma}{T} \boldsymbol{\iota} \boldsymbol{\iota}^\top \right) \mathbf{X} \mathbf{M}_{\mathbf{A}} \right]$

From this comparison we are able to see that both the modified PCA and conventional PCA share a similar functional form of objective function and once $\gamma = -1$, the objective function of modified PCA reduces to objective function of conventional PCA. Moreover, it is obvious that the application of PCA to $\mathbf{X}^\top \left(\mathbf{I}_T + \frac{\gamma}{T} \boldsymbol{\iota} \boldsymbol{\iota}^\top \right) \mathbf{X}$ is equivalent to the application of PCA to $\frac{1}{T} \mathbf{X}^\top \mathbf{X} + \gamma \bar{\mathbf{X}} \bar{\mathbf{X}}^\top$, where $\bar{\mathbf{X}} = \frac{1}{T} \mathbf{X}^\top \boldsymbol{\iota}$.

¹⁰ Actually (88) corresponds to a two-step optimization: for a fixed \mathbf{A} , \mathbf{F} is chosen following the OLS-style first order condition, which implies that $\frac{1}{T} \mathbf{F}^\top \boldsymbol{\iota} = \left(\mathbf{A}^\top \mathbf{A} \right)^{-1} \mathbf{A}^\top \left(\frac{1}{T} \mathbf{X}^\top \boldsymbol{\iota} \right)$, then substituting back \mathbf{F} of this form yields the R.H.S. of (88).

7 Empirical Application

7.1 Data

Nowadays there are several benchmark datasets corresponding to the firm-level characteristics and returns collected at monthly-frequency, including the work done by Green et al. (2017), Gu et al. (2019), Kozak et al. (2020) and Chen and Zimmermann (2020a) (cited as CZ20). Among which CZ20 is in particular by far the most recent and comprehensive one that has successfully covered almost all the major documented anomalies in literature. ¹¹

The major reason why these firm-level characteristics are important is that essentially the discrepancy of stock returns at cross-sectional dimension is attributed to these documented characteristics and actually the expected return is a function of specific form of these characteristics. Conceptionally it is possible for us to regard each characteristic as one measure to distinguish different firms, and for dimension-reduction concern a specific factor structure summarizing the major source of cross-sectional variation is imposed in literature. Factor model or the factor-structure imposed in research is not simply confined to the discussion of stock market but applies broadly in other financial market like mutual funds (Warther, 1995; Barber, Huang, and Odean, 2016; Edelen and Warner, 2001; Berk and Green, 2004; Bergstresser and Poterba, 2004; Del Guercio and Tkac, 2002; Frazzini and Lamont, 2008; Ben-Rephael, Kandel, and Wohl, 2011; Ferson and Kim, 2012; Kogan and Papanikolaou, 2013; Dou, Kogan, and Wu, 2020; Lou, 2012).

For decades after the factor framework proposed by Fama and French (1993) in finance literature, which is broadly acknowledged as the initial attempt to complement CAPM, there has a vast majority of researches with focus on searching new anomalies that lead to abnormal returns that cannot be explained benchmark model. However there is still ongoing debate about how the results claimed in academic research generate impact on the associated anomalies-based investment strategies and despite such a kind of research framework is overwhelming in empirical asset pricing literature, recently it has received widely known criticism either for the underlying methodologies or the credibility of the widely documented empirical results. The initially forwarded question about this issue is what emphasized in Cochrane (2017) that the presence of a vast collection of noisy and highly correlated predictors motivates the adoption of new methodologies instead of the conventional cross-sectional regressions and portfolio sorts. However, the prevalent methodology implemented in finance literature to uncover intervention (publication) effect is based on the direct comparison between the average returns of anomaly-based portfolios before and after publication, rarely is there literature discussing the counterfactual effect, i.e. the return if the anomaly had not been published. By contrast, counterfactual effect is of broad interest in Microeconometrics researches and it is intrinsically deserve more attention and recently there is an evolution accommodating this idea in finance literature like the work done in Pelger and Xiong (2020)

¹¹We acknowledge the codes and data kindly shared by the authors and their efforts on constantly maintaining and updating the data. Both the codes and data are available at the authors' maintained website <https://sites.google.com/site/chenandrewy/open-source-ap?authuser=0>.

For the dataset constructed in CZ20, there are in total 210 portfolios based on different anomalies (characteristics). Specifically, each anomaly is based on a firm-specific variable (characteristic), e.g. the size and book-to-market ratio and then all the stocks traded on U.S. market are sorted into five quintile portfolios based on the corresponding firm-specific characteristic. Return associated with anomaly is Long-short portfolios that buy the highest quintile and sell the lowest quintile portfolio.

One of their claimed result for the simple case when there is only a single factor is as following

$$\sqrt{T}(\tilde{\lambda}_i - \lambda_i) \xrightarrow{d} \begin{cases} \mathcal{N}\left(0, \frac{T}{T_0} \frac{\sigma_e^2}{\sigma_F^2} + 2 \frac{T - T_0}{T_0}\right) & \text{for } i = 1, \dots, N_0 \\ \mathcal{N}\left(0, \frac{\sigma_e^2}{\sigma_F^2}\right) & \text{for } i = N_0 + 1, \dots, N \end{cases}$$

One of the main objective for Pelger and Xiong (2020) is to uncover the difference between common component of units after treatment adoption C_{it}^{treat} and the common component of the synthetic control C_{it}^{ctrl} , that is

$$\tau_{it} = C_{it}^{\text{treat}} - C_{it}^{\text{ctrl}} \tag{91}$$

Generally there are three kinds of patterns associated with data structure with missing values: (i) random missing pattern where whether the entry of data is observed or not does not depend on the entries or other observable covariates; (ii) All the data for treated panel are not observed simultaneously after a specific period; (iii) The periods after which the data for units collected in treated panel are not observed is staggered. This categorization is demonstrated visually as following

[Place [Figure 5](#) about here]

For the inspiring insights on the practical empirical application, Let us first visualize the cumulative return associated with these long-short portfolios based sorting on anomalies (characteristics)

[Place [Figure 6](#) about here]

This may not be that obvious, another direct measure would be calculating and comparing the average returns of different anomaly-based portfolios before and after the publication of specific anomaly.

[Place [Figure 7](#) about here]

Figures and Tables

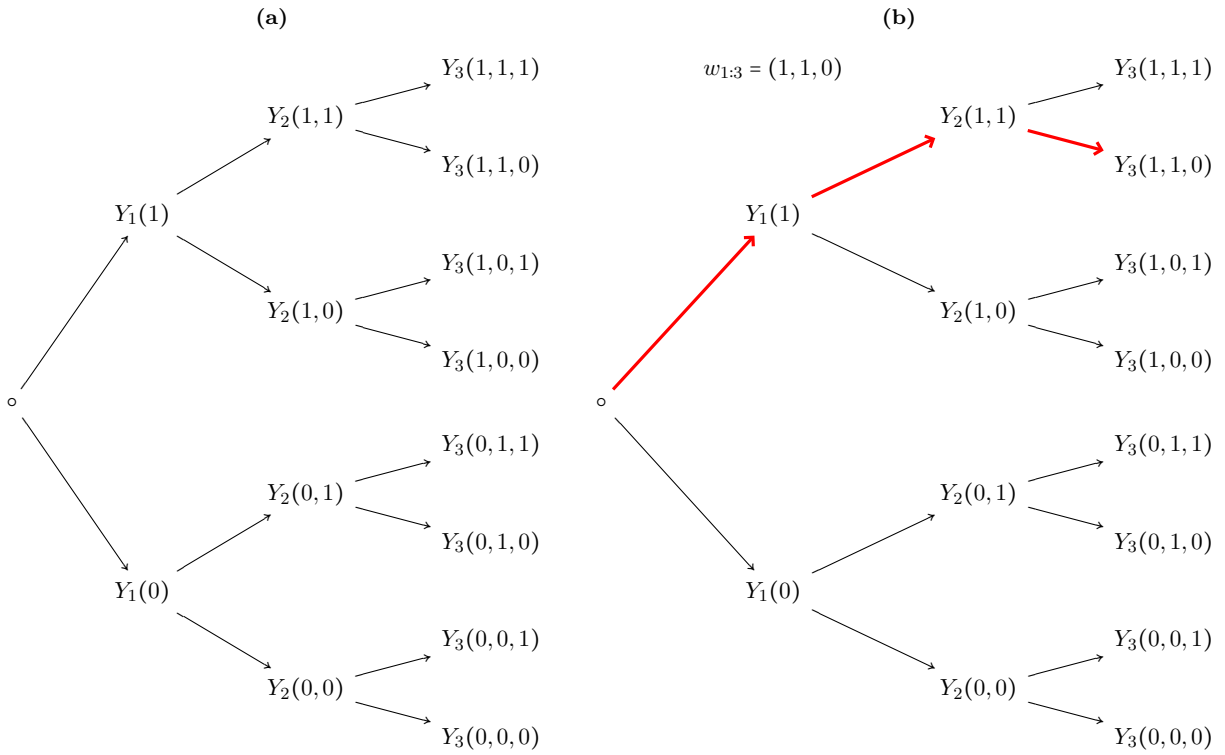


Figure 1. (a) demonstrates all the potential outcomes from intervention for $T = 3$. (b) demonstrates the observed outcome path (possibly not affected by intervention at some period) $Y_{1:3}(w_{1:3})$, indicated by the thick red line.

Table 1. Summary of basic concepts associated variational bayes method

Name	Definition
Variational log likelihood	$M_n(\boldsymbol{\theta}; \mathbf{x}) := \sup_{q(\mathbf{z}) \in \Omega^n} F(q, \boldsymbol{\theta}) = \sup_{q(\mathbf{z}) \in \Omega^n} \int q(\mathbf{z}) \ln \left(\frac{p(\mathbf{x}, \mathbf{z} \boldsymbol{\theta})}{q(\mathbf{z})} \right) d\mathbf{z}$
Variational frequentist estimate (VFE)	$\operatorname{argmax}_{\boldsymbol{\theta}} M_n(\boldsymbol{\theta}, \mathbf{x})$
VB ideal	$\pi^*(\boldsymbol{\theta} \mathbf{x}) = \frac{p(\boldsymbol{\theta}) M_n(\boldsymbol{\theta}; \mathbf{x})}{\int p(\boldsymbol{\theta}) M_n(\boldsymbol{\theta}; \mathbf{x}) d\boldsymbol{\theta}}$
Evidence Lower Bound (ELBO)	$\text{ELBO}(q(\boldsymbol{\theta}, \mathbf{z})) := \iint q(\boldsymbol{\theta}) q(\mathbf{z}) \ln \frac{p(\mathbf{x}, \mathbf{z}, \boldsymbol{\theta})}{q(\boldsymbol{\theta}) q(\mathbf{z})} d\boldsymbol{\theta} d\mathbf{z}$
VB posterior	$q^*(\boldsymbol{\theta}) := \operatorname{argmax}_{q(\boldsymbol{\theta}) \in \Omega^d} \sup_{q(\mathbf{z}) \in \Omega^n} \text{ELBO}(q(\boldsymbol{\theta}, \mathbf{z}))$
VB estimate (VBE)	$\hat{\boldsymbol{\theta}}_n^* = \int \boldsymbol{\theta} \cdot q^*(\boldsymbol{\theta}) d\boldsymbol{\theta}$

Note: (i) $M_n(\boldsymbol{\theta}; \mathbf{x})$ is introduced with subscript n to emphasize the dependency on sample size. (ii) Ω^n is introduced with superscript n to emphasize that the optimization is implemented over the family of distributions of local latent variables and all the alternatives contained in the family of these distributions are factorizable.

Figure 2. In this figure we demonstrate how the posterior probabilities assigned to the first 20 covariates evolve over time. The length of time-series observations is $T = 200$ and the total number of covariates is $p = 200$.

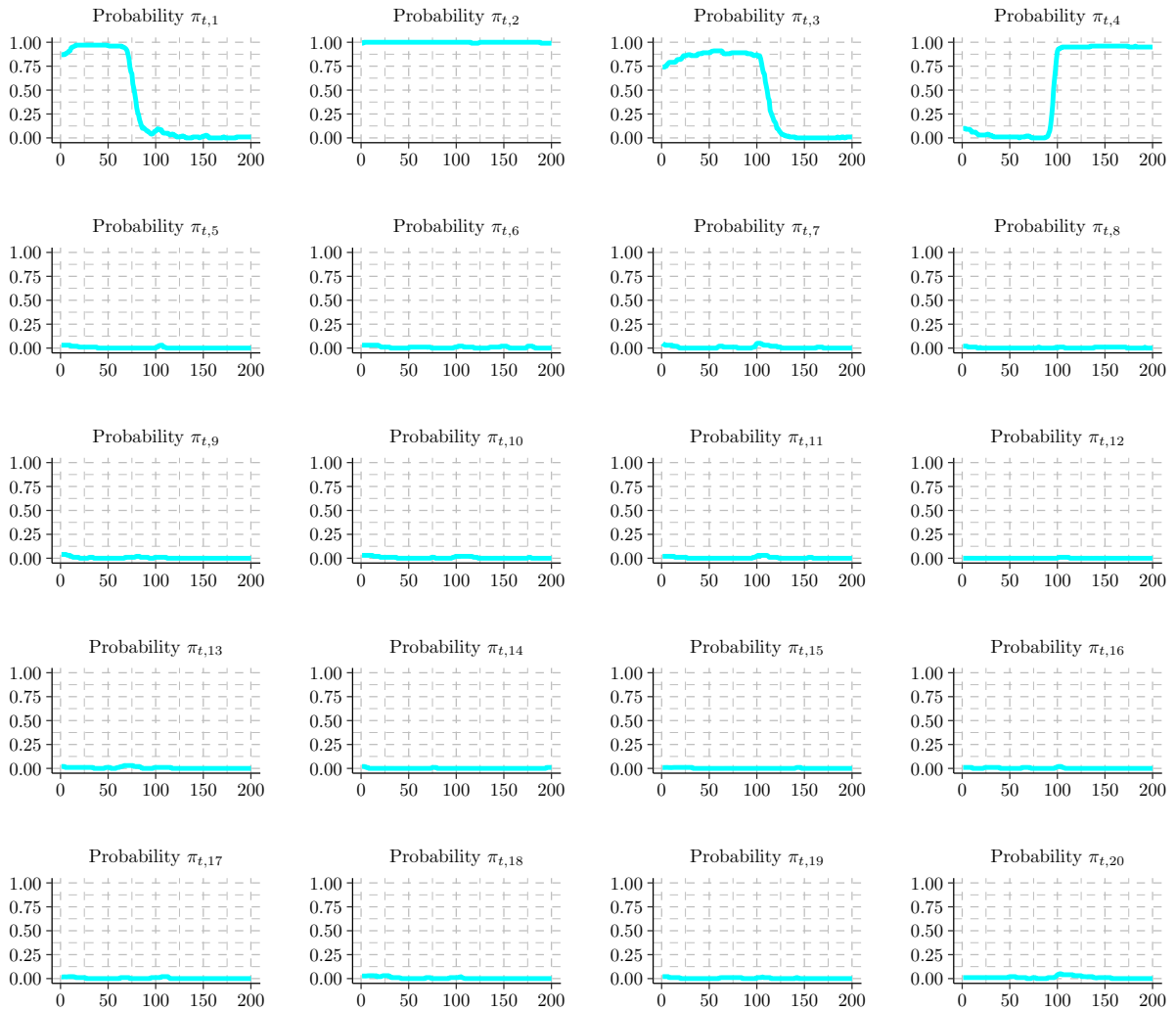


Figure 3. In this figure, we demonstrate how coefficients are estimated (medians over 100 Monte Carlo iterations) over time and the corresponding true values used for generating data as well. Cyan solid lines refer to coefficients fitted from the variational bayes algorithm and black long-dashed lines refer to true values used to generate data. The sub-figure contained in each panel represents how the related values evolve over time. And grey areas refer to the 84% and 16% quantiles (over 100 Monte Carlo iterations) respectively. The length of time-series observations is $T = 200$ and the total number of covariates is $p = 200$.

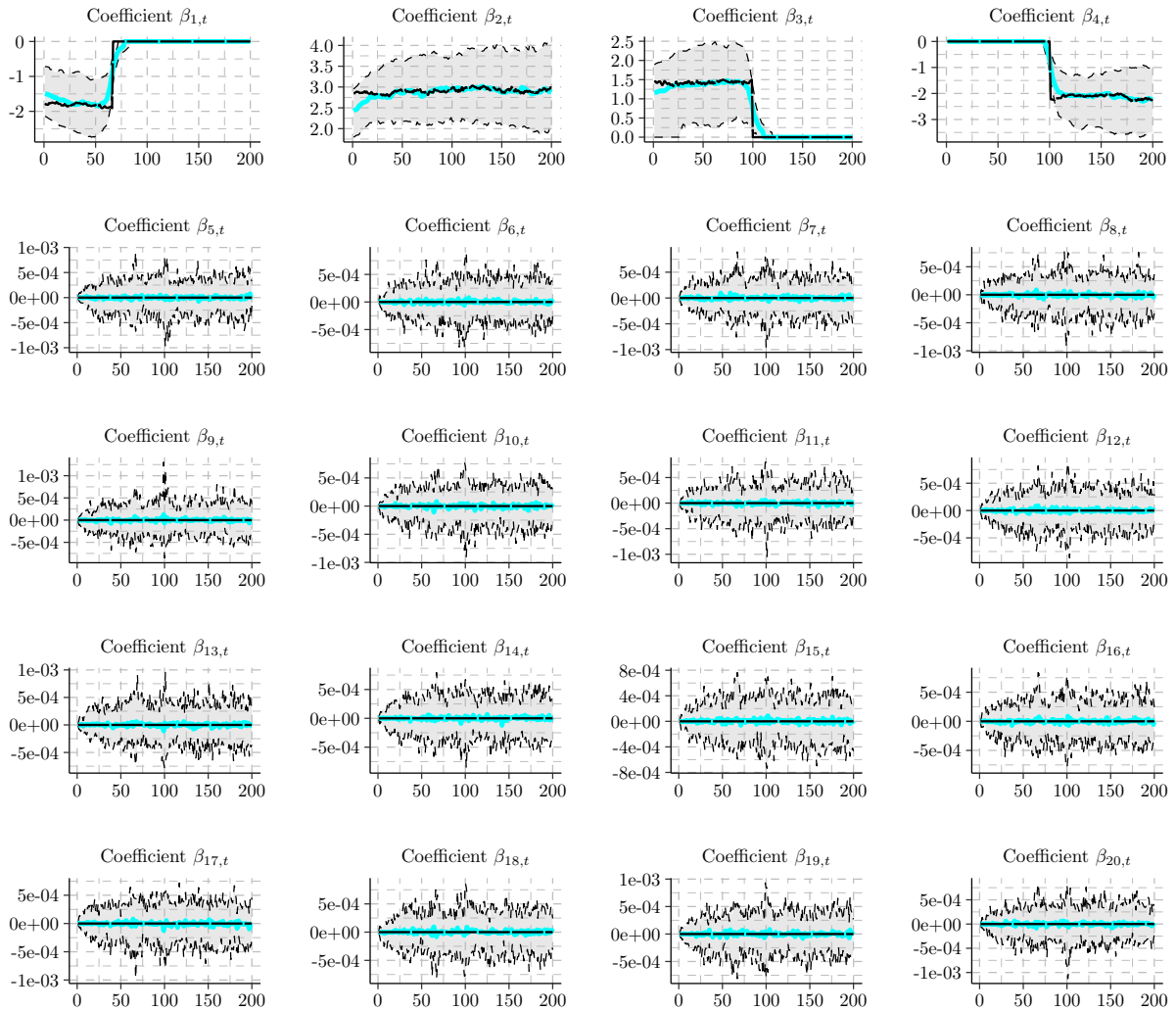


Figure 4. In this figure we demonstrate how the posterior probabilities assigned to the first 20 covariates evolve over time. The length of time-series observations is $T = 200$ and the total number of covariates is $p = 200$.

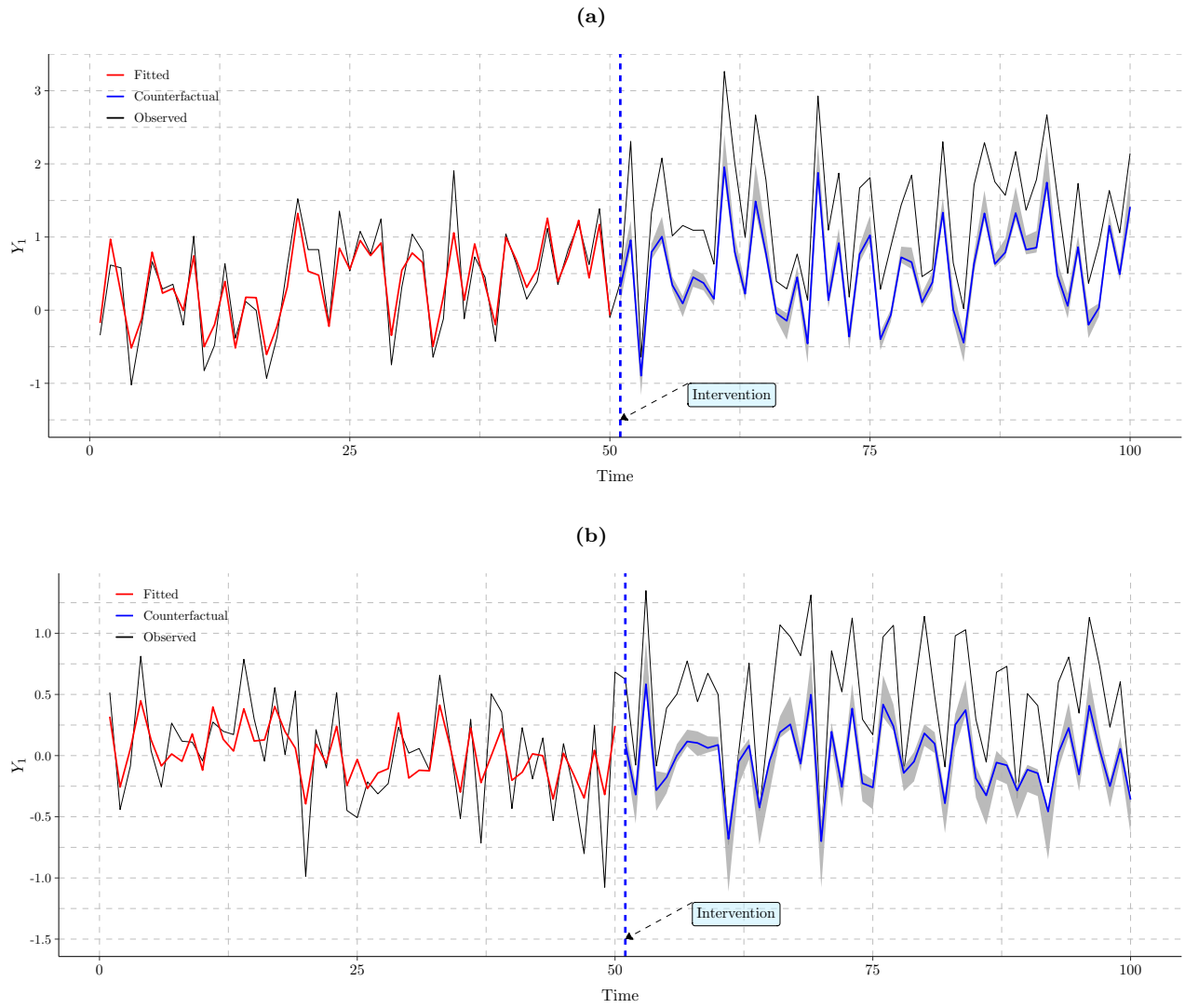


Figure 5. This figure is for demonstrating visually the ideas of three data missing patterns mentioned in the context. Each panel refers to a specific missing pattern, i.e., (a), (b) and (c) refers to random missing pattern, simultaneously missing pattern and staggered missing pattern respectively. To make it as a comparable discussion in panel causal inference setting, horizontal axis refers to the time while the vertical axis refers to the units. Blue shaded entries indicate the missing observations.

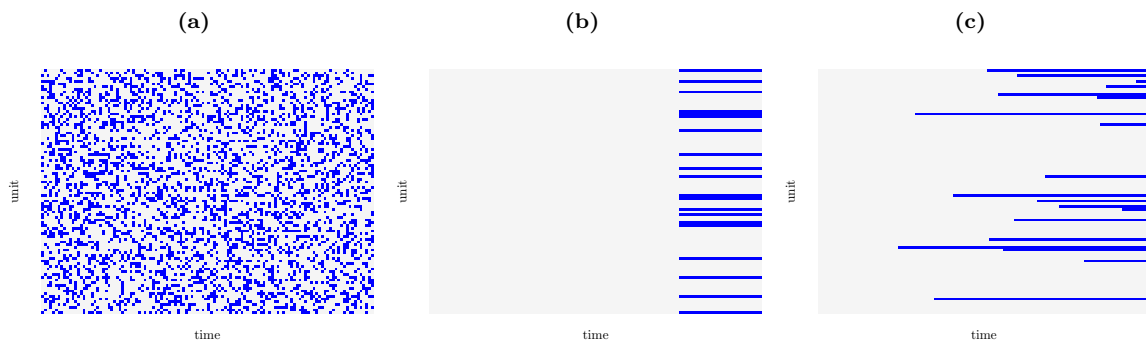
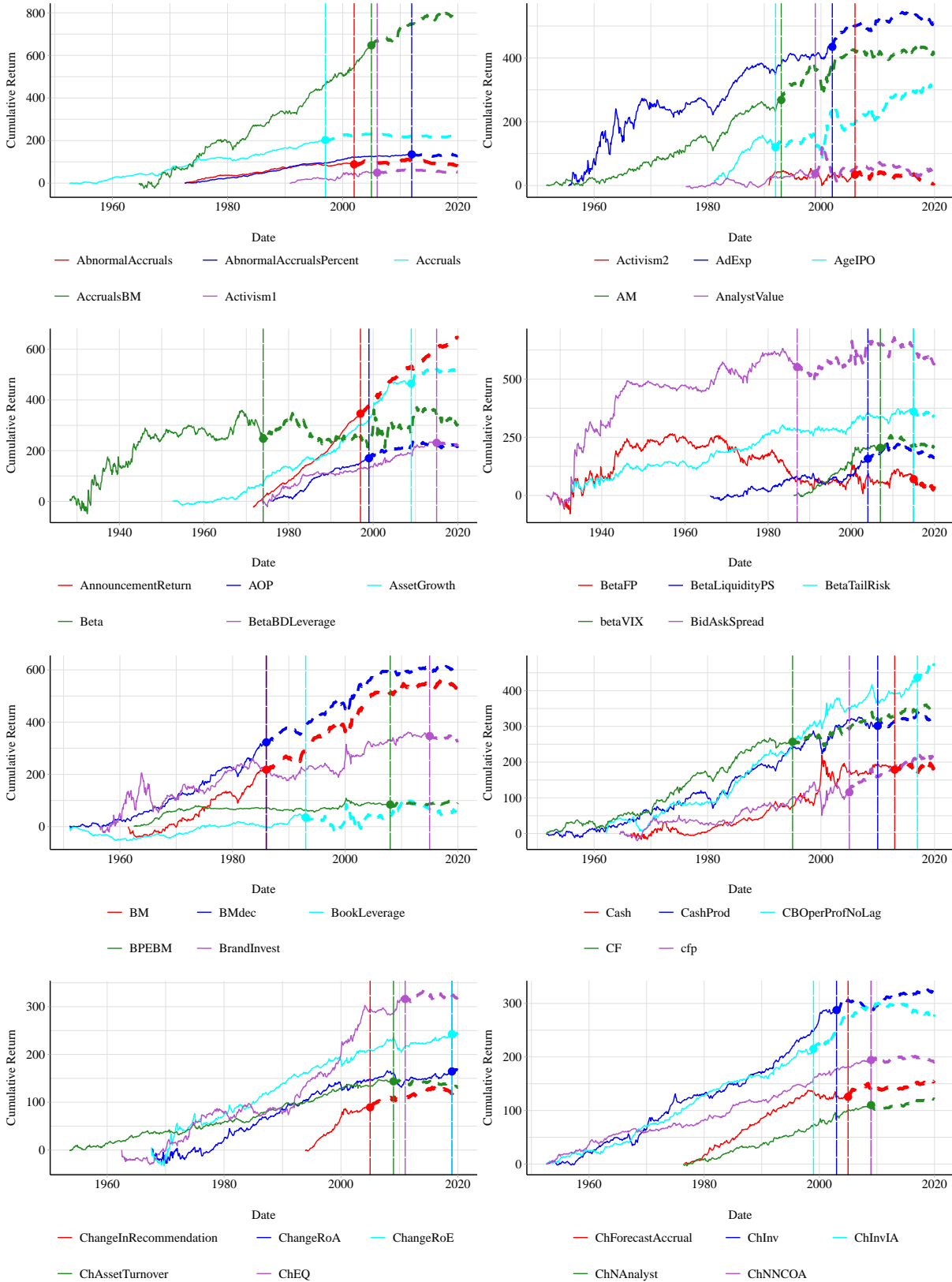
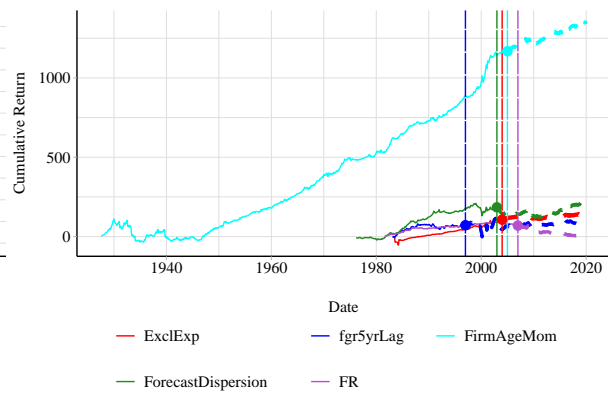
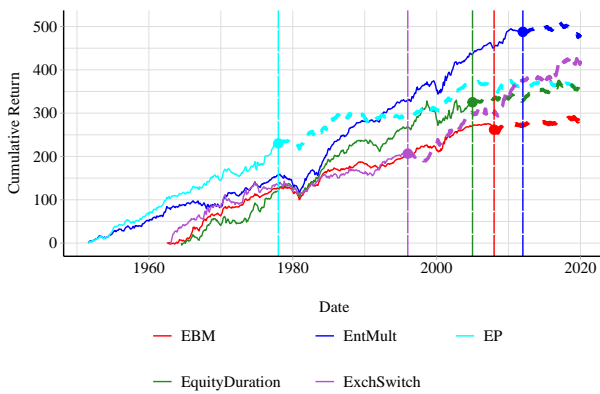
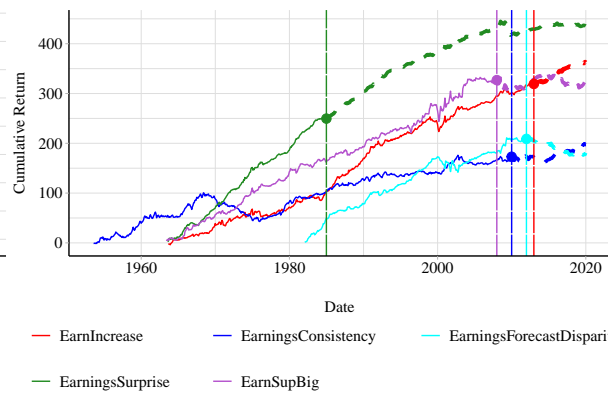
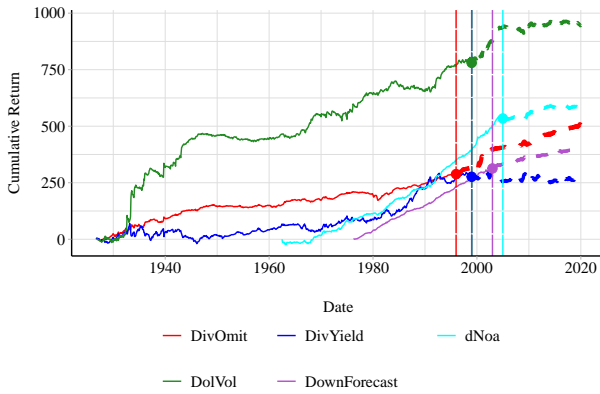
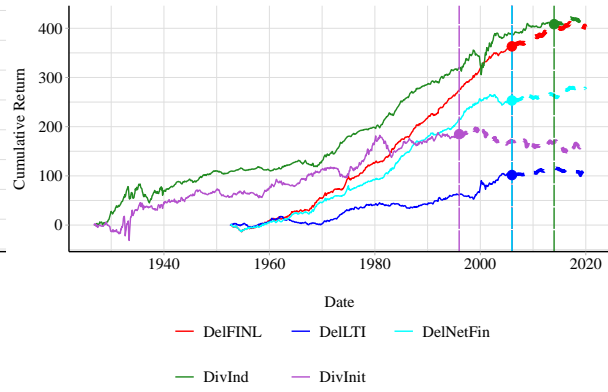
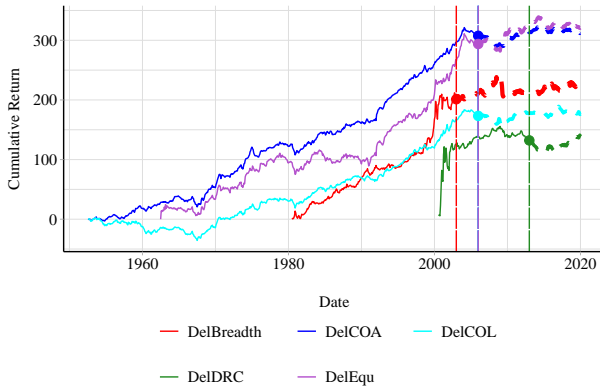
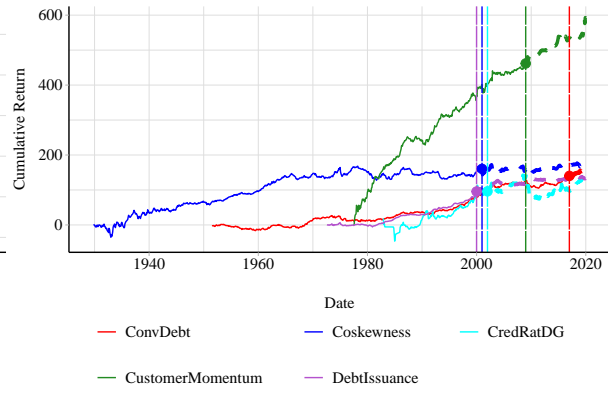
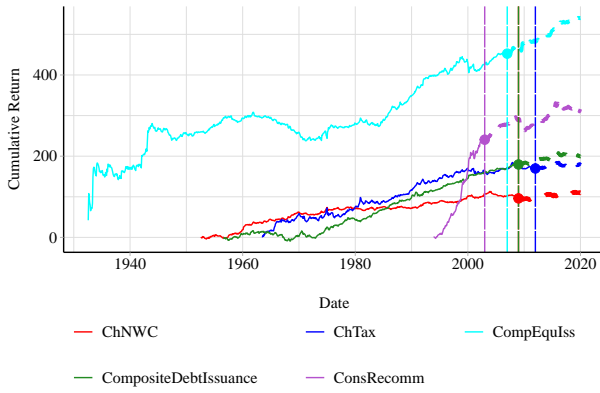
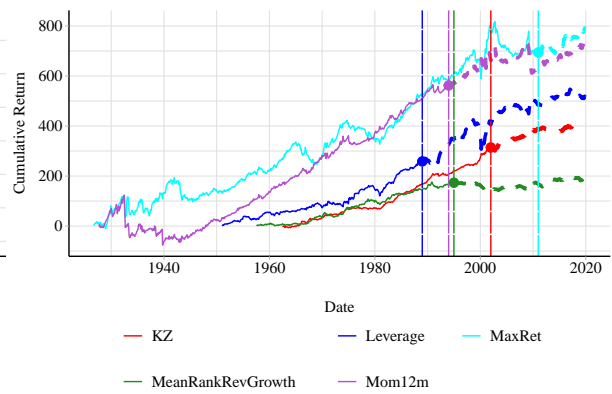
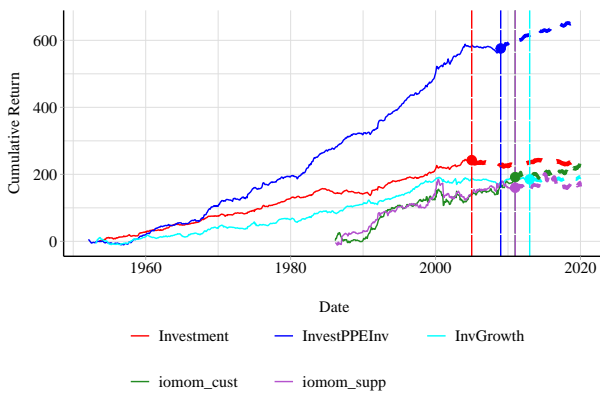
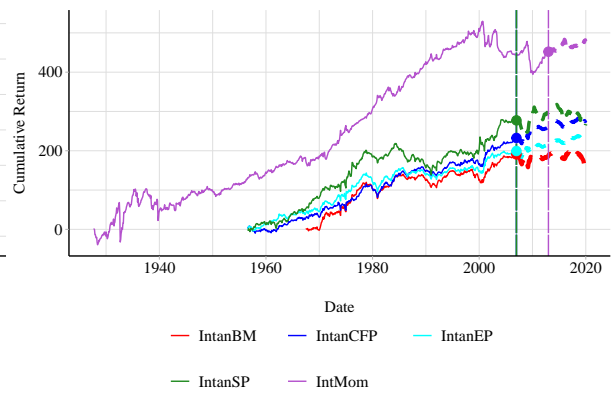
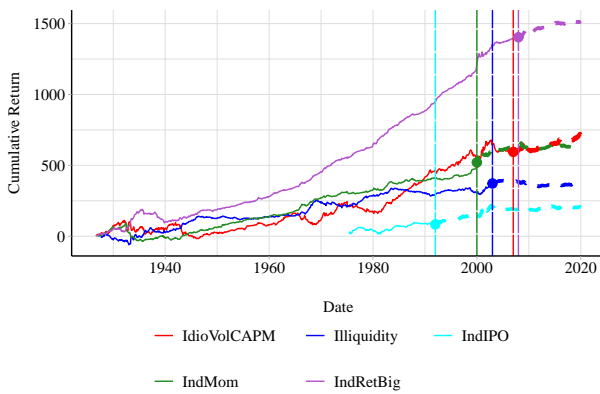
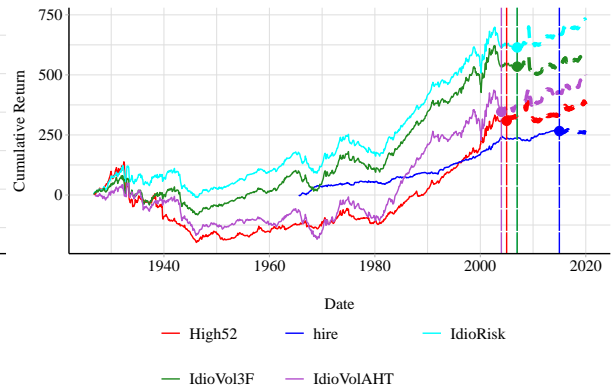
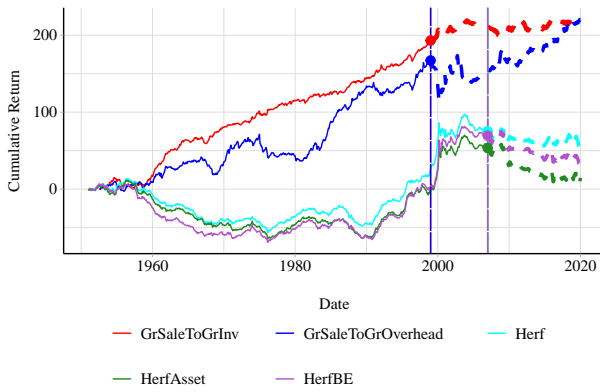
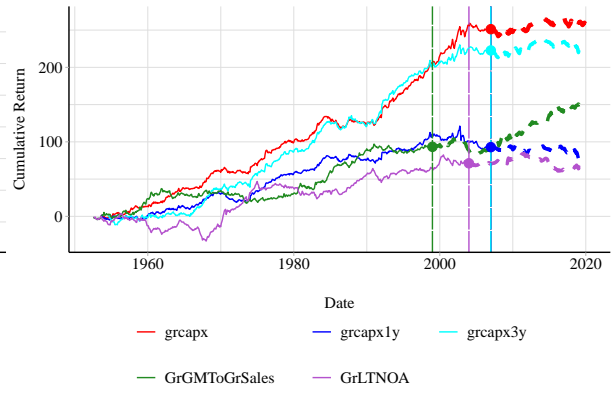
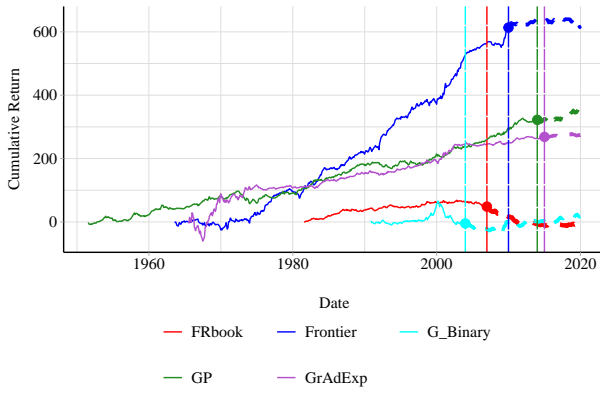
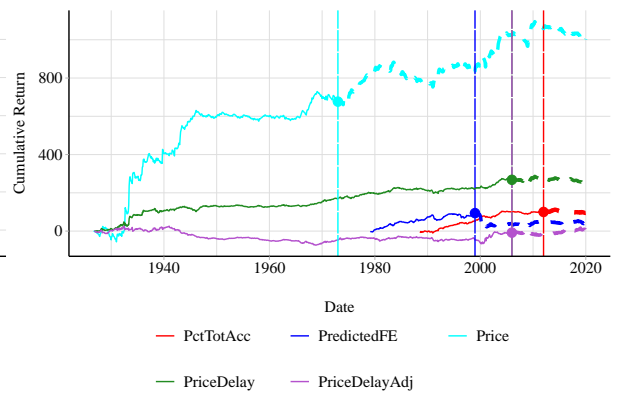
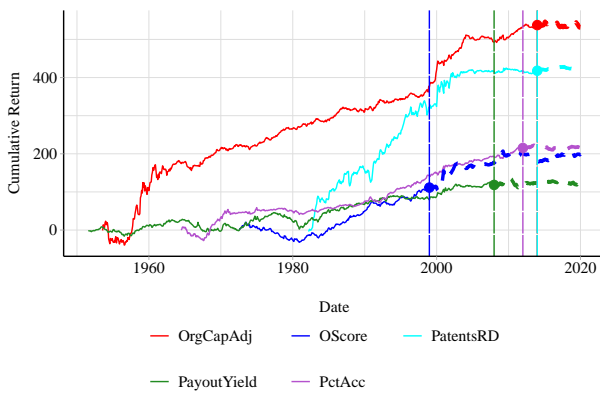
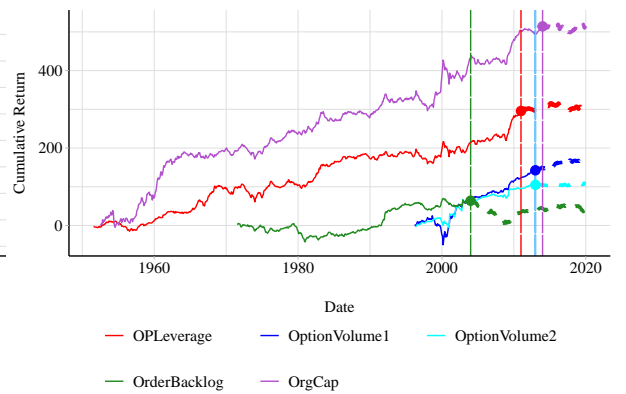
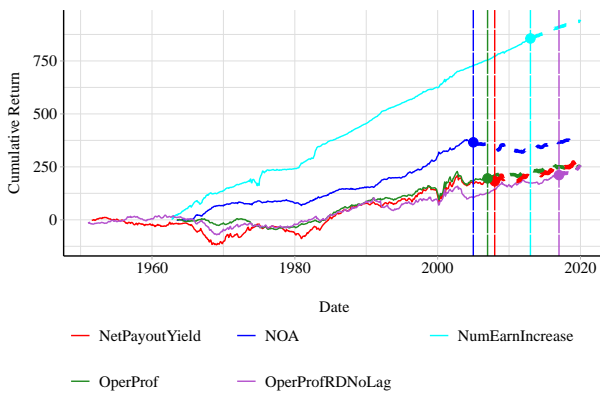
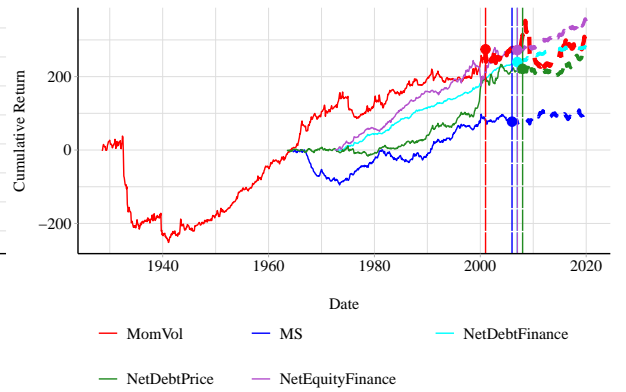
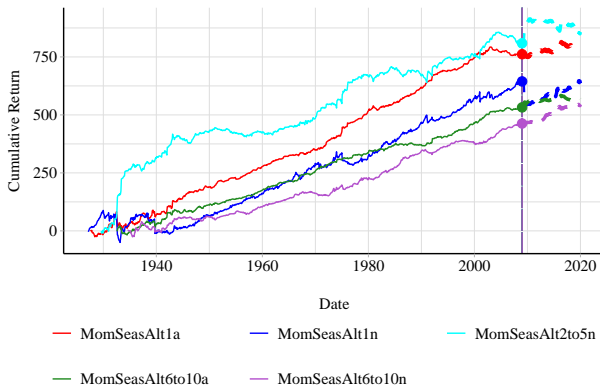
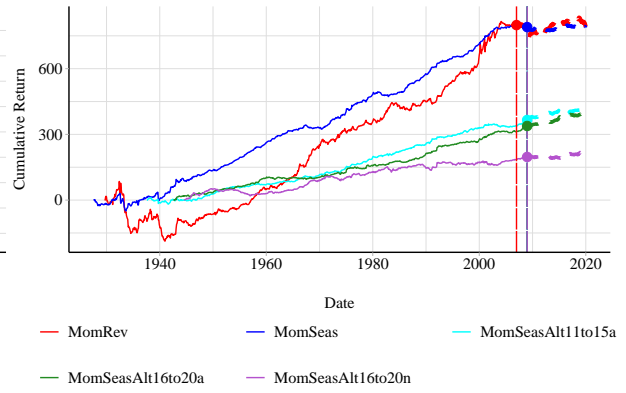
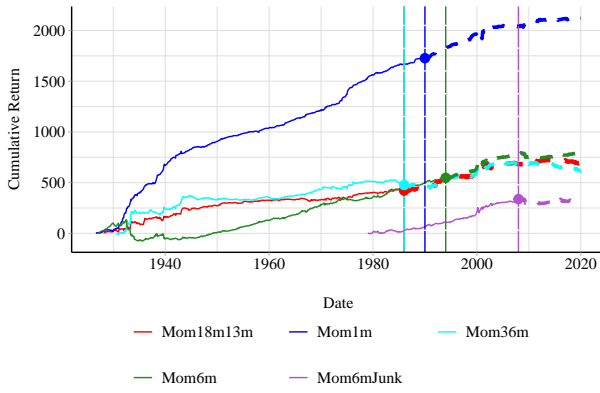


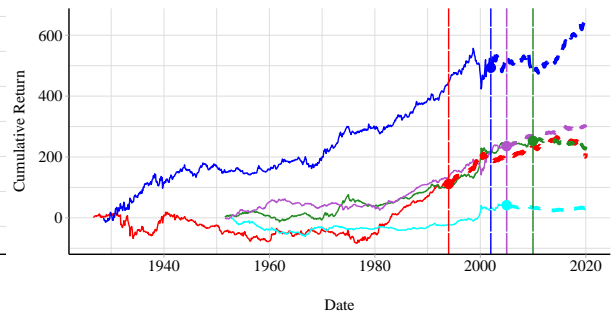
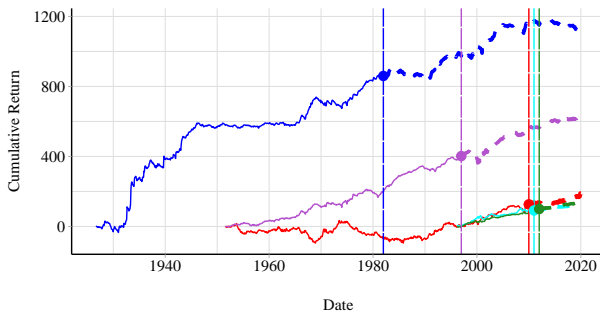
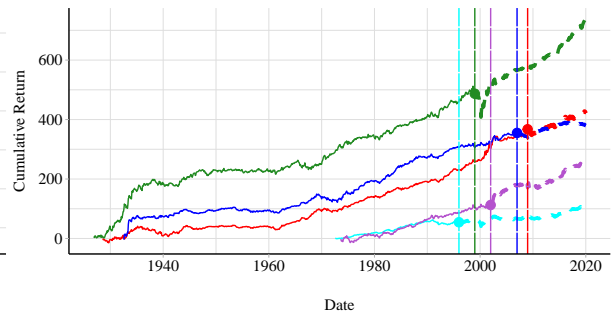
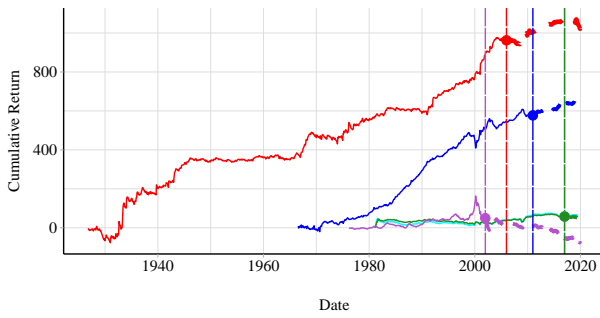
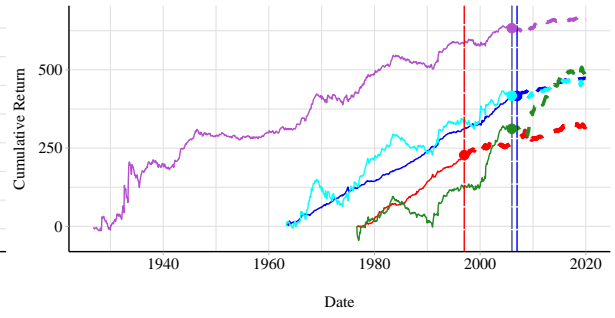
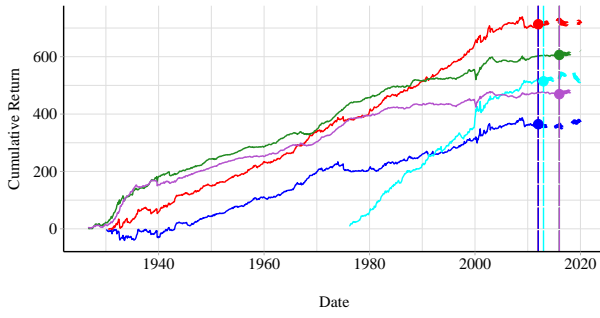
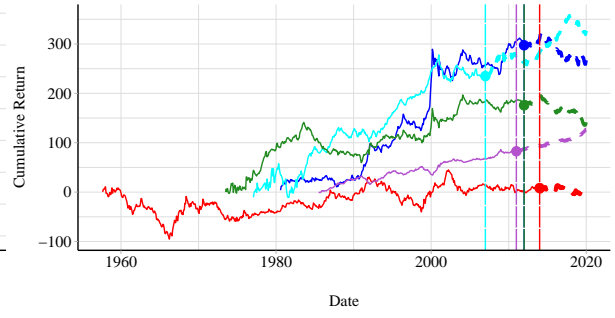
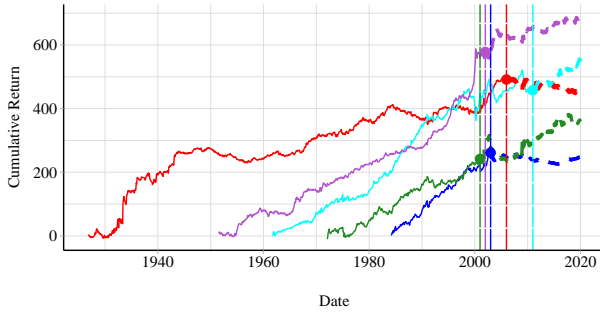
Figure 6

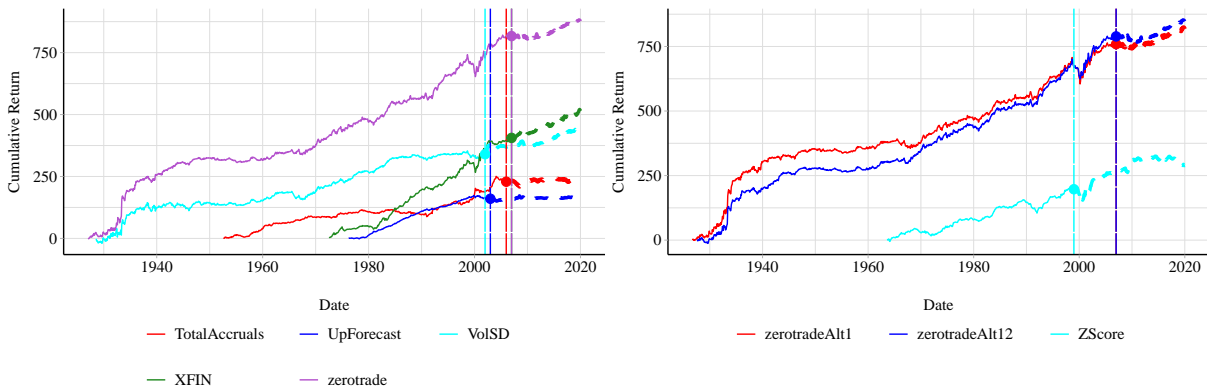






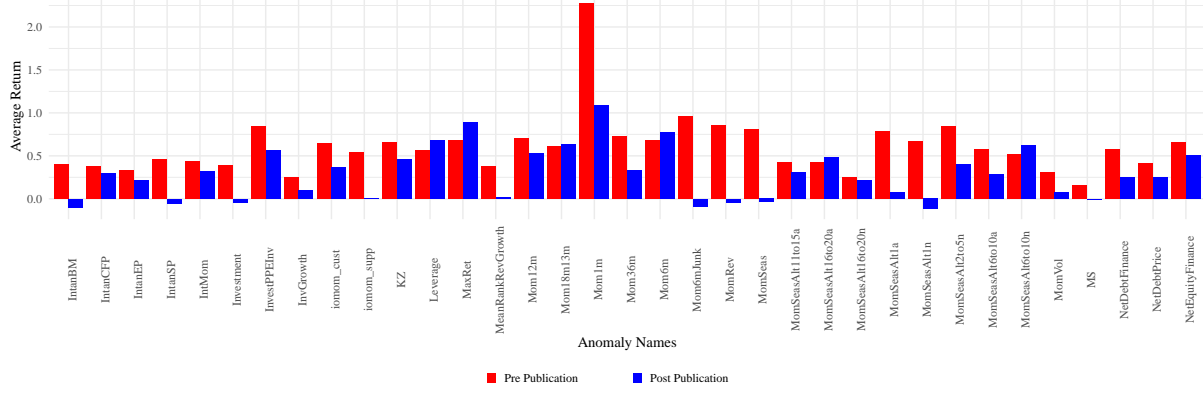
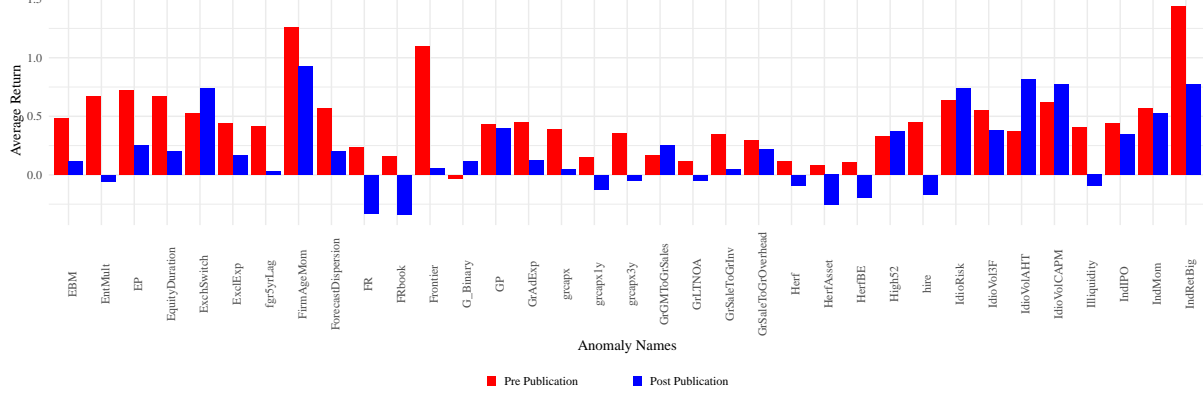
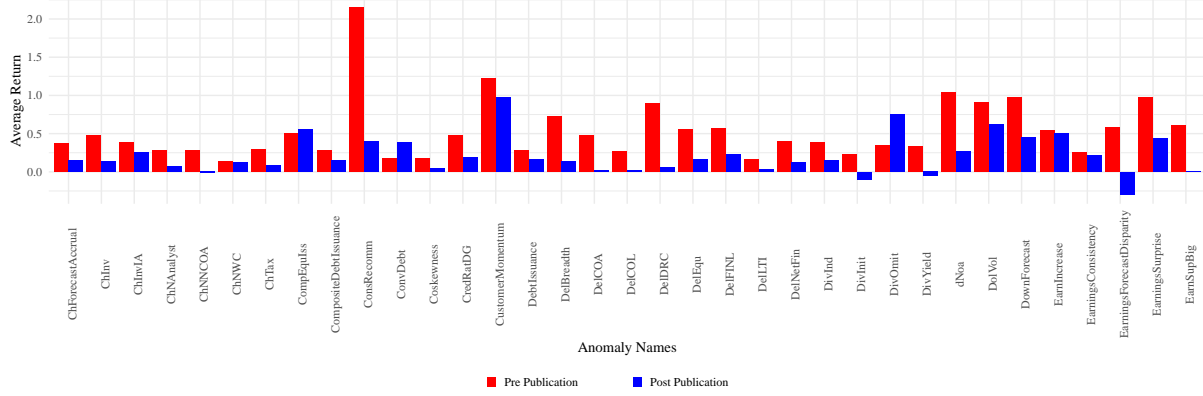
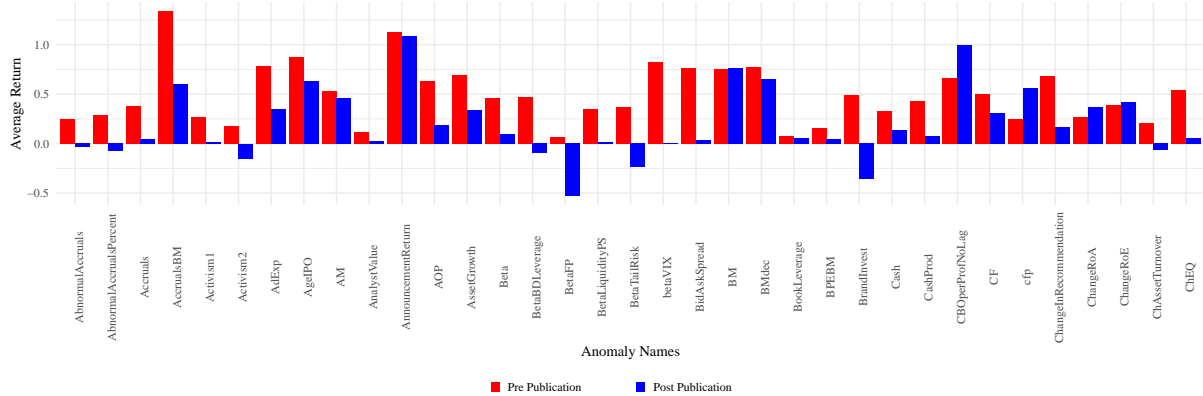


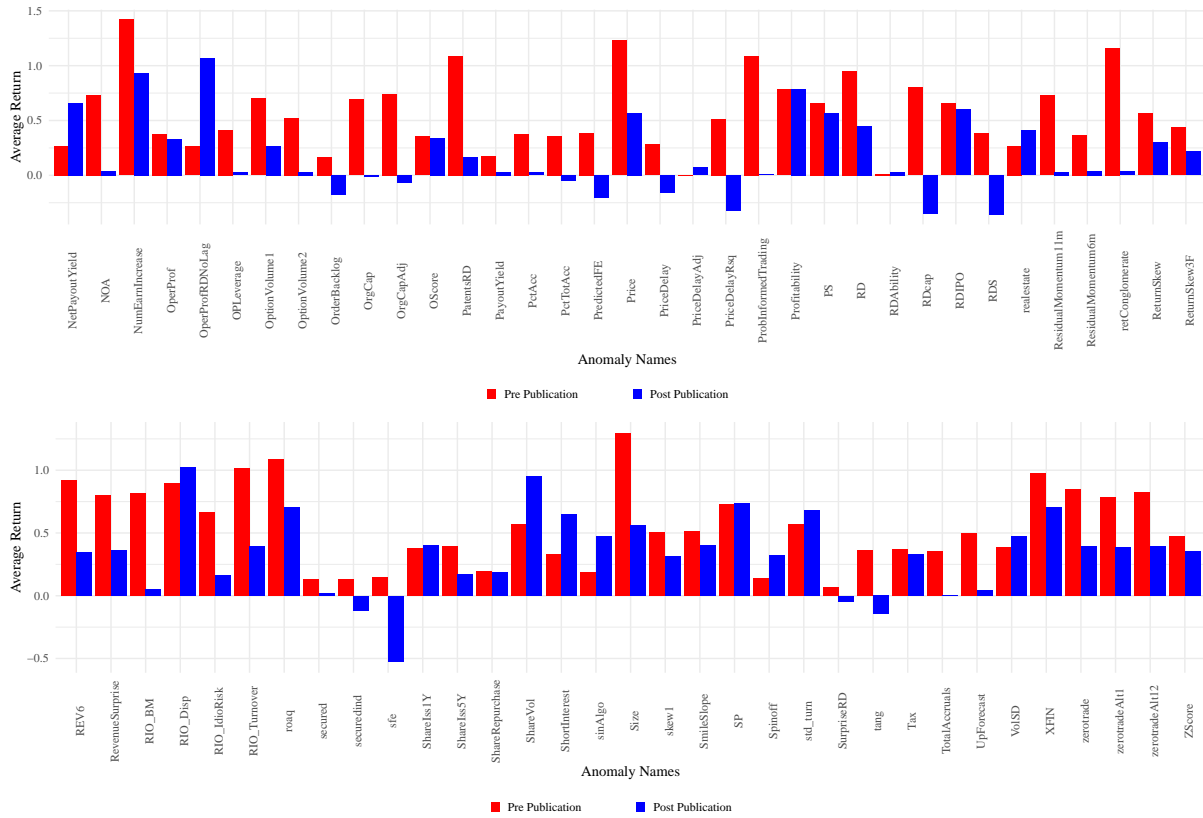




Note: In each panel of the above figure, there are 5 plots of cumulative returns corresponding to different long-short portfolios based on different anomalies (characteristics) respectively. The intersection of vertical line and cumulative return plot indicates timing when the corresponding anomaly was published. Consequently, each cumulative return plot is separated into two parts: the solid part refers to the period before publication of specific anomaly while the dashed part refers to the period after publication of specific anomaly.

Figure 7





Note: In each panel of the above figure, average returns of different anomaly-based long-short portfolios are demonstrated separately. Specifically, the red bar plot indicates the the average return of a specific anomaly-based long-short portfolio before that anomaly was published while the blue bar indicates the average return of long-short portfolio constructed from that anomaly after publication.

References

- ABADIE, A. (2020): “Using Synthetic Controls: Feasibility, Data Requirements, and Methodological Aspects,” Working paper, Massachusetts Institute of Technology, forthcoming in *Journal of Economic Literature*. [Cited on pages 2, 5, and 6.]
- ABADIE, A., A. DIAMOND, AND J. HAINMUELLER (2010): “Synthetic Control Methods for Comparative Case Studies: Estimating the Effect of California’s Tobacco Control Program,” *Journal of the American Statistical Association*, 105, 493–505. [Cited on page 2.]
- ABADIE, A. AND J. GARDEAZABAL (2003): “The Economic Costs of Conflict: A Case Study of the Basque Country,” *American Economic Review*, 93, 113–132. [Cited on page 2.]
- ALQUIER, P. AND J. RIDGWAY (2020): “Concentration of Tempered Posteriors and of Their Variational Approximations,” *The Annals of Statistics*, 48, 1475–1497. [Cited on pages 1, 20, 21, and 27.]
- ARMAGAN, A., D. B. DUNSON, AND M. CLYDE (2011): “Generalized Beta Mixtures of Gaussians,” *Advances in neural information processing systems*, 523–531. [Cited on page 33.]
- ARMSTRONG, T. B., M. KOLESÁR, AND M. PLAGBORG-MØLLER (2020): “Robust Empirical Bayes Confidence Intervals,” Working paper, Yale University and Princeton University. [Cited on page 1.]
- ATHEY, S., M. BAYATI, N. DOUDCHENKOU, G. IMBENS, AND K. KHOSRAVI (2020): “Matrix Completion Methods for Causal Panel Data Models,” Working paper, arXiv:1710.10251. [Cited on page 2.]
- BARBER, B. M., X. HUANG, AND T. ODEAN (2016): 29, 2600–2642. [Cited on page 42.]
- BARNDORFF-NIELSEN, O., J. KENT, AND M. SØRENSEN (1982): “Normal Variance-Mean Mixtures and z Distributions,” *International Statistical Institute (ISI)*, 50, 145–159. [Cited on page 34.]
- BEN-REPHAEL, A., S. KANDEL, AND A. WOHL (2011): “The Price Pressure of Aggregate Mutual Fund Flows,” *The Journal of Financial and Quantitative Analysis*, 46, 585–603. [Cited on page 42.]
- BERGSTRESSER, D. AND J. POTERBA (2004): “Do after-tax returns affect mutual fund inflows,” *Journal of Financial Economics*, 63, 381–414. [Cited on page 42.]
- BERK, J. B. AND R. C. GREEN (2004): “Mutual Fund Flows and Performance in Rational Markets,” *Journal of Political Economy*, 112, 1269–1295. [Cited on page 42.]
- BLEI, D. M., A. KUCUKELBIR, AND J. D. MCAULIFFE (2017): “Variational Inference: A Review for Statisticians,” ArXiv:1601.00670. [Cited on page 1.]
- BLEI, D. M., A. Y. NG, AND M. I. JORDAN (2003): “Latent Dirichlet Allocation,” *Journal of Machine Learning Research*, 3, 993–1022. [Cited on page 8.]

- BOJINOV, I. AND N. SHEPARD (2019): “Time Series Experiments and Causal Estimands: Exact Randomization Tests and Trading,” *Journal of the American Statistical Association*, 114, 1665–1682. [Cited on page 7.]
- BRAUN, M. AND J. MCAULIFFE (2010): “Variational Inference for Large-scale Models of Discrete Choice,” *Journal of American Statistical Association*, 105, 324–334. [Cited on page 1.]
- CARLIN, B. P. AND T. A. LOUIS (2000): *Bayes and Empirical Bayes Methods for Data Analysis*. Chapman & Hall/CRC, New York, NY, 2nd edition. [Cited on page 36.]
- CAROLAN, C. A. (2002): “The Least Concave Majorant of the Empirical Distribution Function,” *The Canadian Journal of Statistics*, 30, 317–328. [Cited on page 37.]
- CARVALHO, C. V., R. MASINI, AND M. C. MEDEIROS (2018): *Journal of Econometrics*, 207, 352–380. [Cited on page 2.]
- CATONI, O. (2004): *Statistical Learning Theory and Stochastic Optimization*. Saint-Flour Summer School on Probability Theory 2001 (Jean Picard ed.), Lecture Notes in Mathematics. Springer. [Cited on page 20.]
- (2007): *PAC-Bayesian supervised classification: the thermodynamics of statistical learning*. Institute of Mathematical Statistics Lecture Notes— Monograph Series, 56. Institute of Mathematical Statistics, Beachwood, OH. [Cited on page 20.]
- CHEN, A. Y. AND T. ZIMMERMANN (2020a): “Open Source Cross-Sectional Asset Pricing,” Working paper, https://papers.ssrn.com/sol3/papers.cfm?abstract_id=3604626. [Cited on pages 42 and 43.]
- (2020b): “Publication Bias and the Cross-Section of Stock Returns,” 10, 249–289. [Cited on page 2.]
- CHERNOZHUKOV, V., K. WÜTHRICH, AND Y. ZHU (2020a): “An Exact and Robust Conformal Inference Method for Counterfactual and Synthetic Controls,” Working paper, arXiv:1712.09089. [Cited on page 2.]
- (2020b): “Practical and Robust t-test based Inference for Synthetic Control and Related Methods,” Working paper, arXiv:1812.10820. [Cited on page 2.]
- CHIPMAN, H. A., E. I. GEORGE, AND R. E. MCCULLOCH (1998): “Bayesian CART Model Search,” *Journal of the American Statistical Association*, 93, 935–948. [Cited on page 5.]
- (2010): “BART: Bayesian Additive Regression Trees,” *The Annals of Applied Statistics*, 4, 266–298. [Cited on page 5.]
- COCHRANE, J. H. (2017): “Presidential Address: Discount rates,” *Journal of Finance*, 66, 1047–1108. [Cited on page 42.]

- COOLEY, T. F. AND E. C. PRESCOTT (1976): “Estimation in the Presence of Stochastic Parameter Variation,” *Econometrica*, 44, 167–184. [Cited on page 1.]
- DEL GUERCIO, D. AND P. A. TKAC (2002): “Star Power: The Effect of Morningstar Ratings on Mutual Fund Flow,” *The Journal of Financial and Quantitative Analysis*, 43, 907–936. [Cited on page 42.]
- DEMPSTER, A. P., N. M. LAIRD, AND B. RUBIN (1977): “Maximum Likelihood from Incomplete Data via the EM Algorithm,” *Journal of Royal Statistical Society. Series B (Methodological)*, 39, 1–38. [Cited on pages 10 and 11.]
- DENISON, D. G. T., B. K. MALLICK, AND A. F. M. SMITH (1998): “A Bayesian CART Algorithm,” *Biometrika*, 85, 363–377. [Cited on page 5.]
- DOU, W. W., L. KOGAN, AND W. WU (2020): “Common Fund Flows: Flow Hedging and Factor Pricing,” Working paper, https://papers.ssrn.com/sol3/papers.cfm?abstract_id=3543675. [Cited on page 42.]
- DOUDCHENKO, N. AND G. W. IMBENS (2017): “Balancing, Regression, Difference-In-Differences and Synthetic Control Methods: A Synthesis,” Working paper, Stanford University. [Cited on page 2.]
- EDDELBUETTEL, D. (2013): *Seamless R and C++ Integration with Rcpp*. Use R! Series. Springer. [Cited on page 29.]
- EDELEN, R. M. AND J. B. WARNER (2001): “Aggregate Price Effects of Institutional Trading: A Study of Mutual Fund Flow and Market Returns,” *Journal of Financial Economics*, 59, 195–220. [Cited on page 42.]
- ENGLE, R., S. GIGLIO, H. LEE, B. KELLY, AND J. STROEBEL (2020): “Hedging Climate Change News,” *Review of Financial Studies*, 33, 1184–1216. [Cited on page 24.]
- FAMA, E. F. AND K. R. FRENCH (1993): “Common Risk Factors in the Returns on Stocks and Bonds,” *Journal of Financial Economics*, 33, 3 – 56. [Cited on pages 10 and 42.]
- (1996): “Multifactor Explanations of Asset Pricing Anomalies,” *The Journal of Finance*, 51, 55–84. [Cited on page 10.]
- (2015): “A Five-factor Asset Pricing Model,” *Journal of Financial Economics*, 116, 1 – 22. [Cited on page 10.]
- FERSON, W. E. AND M. S. KIM (2012): “The factor structure of mutual fund flows,” *International Journal of Portfolio Analysis and Management*, 1, 112–143. [Cited on page 42.]
- FRAZZINI, A. AND O. A. LAMONT (2008): “Dumb money: Mutual fund flows and the cross-section of stock returns,” *Journal of Financial Economics*, 88, 299–322. [Cited on page 42.]

- GENTZKOW, M., B. KELLY, AND M. TADY (2019): “Text as Data,” *Journal of Economic Literature*, 57, 535–574. [Cited on page 23.]
- GEORGE, E. I. AND R. E. MCCULLOCH (1993): “Variable Selection via Gibbs Sampling,” *Journal of the American Statistical Association*, 88, 881–889. [Cited on pages 15 and 28.]
- (1997): “Approaches for Bayesian Variable Selection,” *Journal of the American Statistical Association*, 7, 339–373. [Cited on pages 15, 16, and 17.]
- GREEN, J., J. R. M. HAND, AND X. F. ZHANG (2017): “The Characteristics that Provide Independent Information about Average U.S. Monthly Stock Returns,” *The Review of Financial Studies*, 30, 4389–4436. [Cited on page 42.]
- GRIFFITHS, T. L. AND M. STEYVERS (2004): “Finding Scientific Topics,” *Proceedings of the National Academy of Sciences of the United States of America*, 101, 5228–5235. [Cited on page 10.]
- GU, S., B. KELLY, AND D. XIU (2019): “Empirical Asset Pricing via Matching Learning,” Working paper, forthcoming in the *The Review of Financial Studies*. [Cited on page 42.]
- GUO, F., X. WANG, K. FAN, T. BRODERICK, AND D. B. DUNSON (2016): “Boosting Variational Inference,” ArXiv:1611.05559. [Cited on page 1.]
- HANSEN, B. E. (2016): “Efficient Shrinkage in Parametric Models,” *Journal of Econometrics*, 190, 115–132. [Cited on page 1.]
- JORDAN, M. I., Z. GHAHRAMANI, T. S. JAAKKOLA, AND L. K. SAUL (1999): “An Introduction to Variational Methods for Graphical Models,” *Machine Learning*, 37, 183–233. [Cited on page 1.]
- KOGAN, L. AND D. PAPANIKOLAOU (2013): “Firm Characteristics and Stock Returns: The Role of Investment-Specific Shocks,” *The Review of Financial Studies*, 26, 2718–2759. [Cited on page 42.]
- KOOP, G. AND D. KOROBILIS (2020): “Bayesian Dynamic Variable Selection in High Dimensions,” R&R in *Journal of Econometrics*. [Cited on pages 5, 27, and 29.]
- KOWAL, D. R., D. S. MATTESON, AND D. RUPPERT (2019): “Dynamic Shrinkage Processes,” *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 81, 781–804. [Cited on pages 27 and 31.]
- KOZAK, S., S. NAGEL, AND S. SANTOSH (2020): “Shrinking the Cross Section,” *Journal of Financial Economics*, 135, 271–292. [Cited on page 42.]
- LETTU, M. AND M. PELGER (2020a): “Estimating Latent Asset-Pricing Factors,” Forthcoming in *The Journal of Finance*. [Cited on pages 38 and 41.]
- (2020b): “Factors that Fit the Time-Series and Cross-Section of Stock Returns,” *Review of Financial Studies*, 33, 2274–2325. [Cited on page 38.]

- LOU, D. (2012): “A flow-based explanation for return predictability,” *The Review of Financial Studies*, 25, 3457–3489. [Cited on page 42.]
- MCLACHLAN, G. J. AND T. KRISHNAN (2008): *The EM Algorithms and Extensions*. New Jersey: John Wiley & Sons. [Cited on page 10.]
- MCLEAN, R. D. AND J. PONTIFF (2016): “Does Academic Research Destroy Stock Return Predictability,” *Journal of Finance*, 71. [Cited on page 2.]
- MORRIS, C. N. (1983): “Parametric Empirical Bayes Inference: Theory and Applications,” *Journal of the American Statistical Association*, 78, 47–55. [Cited on pages 1 and 36.]
- PELGER, M. AND R. XIONG (2020): “Large Dimensional Latent Factor Modeling with Missing Observations and Applications to Causal Inference,” Working paper, <https://mpelger.people.stanford.edu/research>. [Cited on pages 2, 42, and 43.]
- PHAN, X., L. NGUYEN, AND S. HORIGUCHI (2008): “Learning to Classify Short and Sparse Text & Web with Hidden Topics from Large-Scale Data Collections,” In *Proceedings of the 17th International World Wide Web Conference (WWW 2008)*, 91–100, Beijing, China. [Cited on page 10.]
- POLSON, N. G., J. G. SCOTT, AND J. WINDLE (2013): “Bayesian Inference for Logistic Models using Pólya-Gamma Latent Variables,” *Journal of the American Statistical Association*, 108, 1339–1349. [Cited on pages 34 and 36.]
- ROBBINS, H. AND S. MONRO (1951): “A Stochastic Approximation Method,” *Annals of Mathematical Statistics*, 22, 400–407. [Cited on page 24.]
- ROČKOVÁ, V. (2019): “On Semi-parametric Bernstein-von Mises Theorems for BART,” Working paper. [Cited on page 5.]
- ROČKOVÁ, V. AND E. I. GEORGE (2014): “EMVS: The EM Approach to Bayesian Variable Selection,” *Journal of the American Statistical Association*, 109, 828–846. [Cited on page 15.]
- ROČKOVÁ, V. AND E. SAHA (2019): “On Theory for BART,” Working paper, 22nd International Conference on Artificial Intelligence and Statistics. [Cited on page 5.]
- ROČKOVÁ, V. AND S. VAN DER PAS (2019): “Posterior Concentration for Bayesian Regression Trees and Forests,” Manuscript, just accepted for *Annals of Statistics*. [Cited on page 5.]
- TRAN, M.-N., D. J. NOTT, AND R. KOHN (2017): “Variational Bayes With Intractable Likelihood,” *Journal of Computational and Graphical Statistics*, 26, 873–882. [Cited on page 29.]
- VAN ERVEN, T. AND P. HARREMOËS (2014): “Rényi Divergence and Kullback-Leibler Divergence,” *IEEE Transactions on Information Theory*, 60, 3797–3820. [Cited on page 19.]

- WAINWRIGHT, M. J. AND M. I. JORDAN (2008): “Graphical Models, Exponential Families, and Variational Inference,” *Foundations and Trends in Machine Learning*, 1, 1–305. [Cited on page 10.]
- WANG, Y. AND D. M. BLEI (2019): “Frequentist Consistency of Variational Bayes,” *Journal of the American Statistical Association*, 114, 1147–1161. [Cited on pages 1 and 27.]
- WARTHER, V. A. (1995): “Aggregate Mutual Fund Flows and Security Returns,” *Journal of Financial Economics*, 39, 209–235. [Cited on page 42.]
- WINN, J. AND M. BISHOP (2005): “Variational Message Passing,” *Journal of Machine Learning Research*, 6, 661–694. [Cited on page 1.]
- ZEILER, M. D. (2012): “ADADELTA: An Adaptive Learning Rate Method,” ArXiv:1212.5701. [Cited on page 25.]